# Supervised learning for computer vision: Kernel methods & sparse methods

**Francis Bach**

*SIERRA Project-team, INRIA - Ecole Normale Supérieure*

CVML Summer school, July 2012

# Machine learning

- **Supervised** learning

  - Predict $y \in \mathcal{Y}$ from $x \in \mathcal{X}$, given observations $(x_i, y_i)$, $i = 1, \ldots, n$
  - *Classification, regression*

- **Unsupervised** learning

  - Find structure in $x \in \mathcal{X}$, given observations $x_i$, $i = 1, \ldots, n$
  - *Clustering, dimension reduction*

- Application to many problems and data types:

  - **Computer vision**
  - Bioinformatics, text processing, audio processing
  - etc.

- Specifity: exchanges between **theory** / **algorithms** / **applications**

# Machine learning for computer vision

- Multiplication of digital media

- Many different tasks to be solved

  – Associated with different machine learning problems
  – Massive data to learn from

- **Machine learning is not limited to binary classification!**

# Image retrieval
$\Rightarrow$ **Classification, ranking, outlier detection**

# Image retrieval
## Classification, ranking, outlier detection

# Image retrieval
## Classification, ranking, outlier detection

# Image annotation
## Classification, clustering

# Object recognition ⇒ Multi-label classification

# Personal photos
## ⇒ Classification, clustering, visualization

# Image recognition beyond Flickr
## Digital historical archives



Monsieur, Vous êtes averti de porter samedi prochain 26 janvier quarante écus dans un trou qui est au pied de la croix Montelay sous peine d'avoir la tête cassée à l'heure que vous y penserez le moins. Si l'on ne vous rencontre point vous êtes assuré que le feu sera mis chez vous. Sil en est parlé à qui que ce soit la tête cassée vous aurez.

Archives du Val dOise - 1737

# Machine learning for computer vision

- Multiplication of digital media

- Many different tasks to be solved

  - Associated with different machine learning problems
  - Massive data to learn from

- Similar situations in many fields (e.g., bioinformatics)

# Machine learning for bioinformatics (e.g., proteins)



Primary protein structure
is sequence of a chain of amino acids

Amino Acids

Amino group
$NH_2$

$H$—$C$—$COOH$

R

R group

Acidic
carboxyl
group

Amino Acid

1. Many learning tasks on proteins

   - Classification into functional or structural classes
   - Prediction of cellular localization and interactions

2. Massive data

# Machine learning for computer vision

- Multiplication of digital media

- Many different tasks to be solved

  - Associated with different machine learning problems
  - Massive data to learn from

- Similar situations in many fields (e.g., bioinformatics)

  $\Rightarrow$ **Machine learning for high-dimensional structured data**

# Why not simply memorizing everything?
## Brute force nearest-neighbor classification

# Why is visual recognition difficult?
## Simple problems

• Few object classes with low variability within classes

# Object recognition in everyday images

# Why is visual recognition difficult?
## Real/complex problems

- **Many object classes with high variability within classes**

# Supervised machine learning from examples
## Why is it difficult/interesting?

- Why not simply memorizing everything?

- Problem of generalization

- **Curse of dimensionality**

  – For data in dimension $p$, without assumptions, $2^p$ observations are needed

- **No free lunch theorem**

  – No algorithm can be the best all the time

- **Prior knowledge is needed**

# Supervised learning and regularization

- Data: $x_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$, $i = 1, \ldots, n$

- Minimize with respect to function $f \in \mathcal{F}$:

$$\sum_{i=1}^{n} \ell(y_i, f(x_i)) \qquad + \qquad \frac{\lambda}{2}\|f\|^2$$

<div style="color:red; text-align:center">Error on data      +      Regularization</div>

<div style="color:blue; text-align:center">Loss & function space ?       Norm ?</div>

- Two theoretical/algorithmic issues:

  1. Loss
  2. Function space / norm

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Losses for regression (Shawe-Taylor and Cristianini, 2004)

- **Response**: $y \in \mathbb{R}$, prediction $\hat{y} = f(x)$,

  - quadratic (square) loss $\ell(y, f(x)) = \frac{1}{2}(y - f(x))^2$
  - Not many reasons to go beyond square loss!

# Losses for regression (Shawe-Taylor and Cristianini, 2004)

- **Response**: $y \in \mathbb{R}$, prediction $\hat{y} = f(x)$,

  - quadratic (square) loss $\ell(y, f(x)) = \frac{1}{2}(y - f(x))^2$
  - Not many reasons to go beyond square loss!

- Other convex losses "with added benefits"

  - $\varepsilon$-insensitive loss $\ell(y, f(x)) = (|y - f(x)| - \varepsilon)_+$
  - Hüber loss (mixed quadratic/linear): robustness to outliers

# Losses for classification (Shawe-Taylor and Cristianini, 2004)

- **Label** : $y \in \{-1, 1\}$, prediction $\hat{y} = \text{sign}(f(x))$

  – loss of the form $\ell(y, f(x)) = \ell(yf(x))$
  – "True" cost: $\ell(yf(x)) = 1_{yf(x)<0}$
  – Usual <span style="color:red">convex</span> costs:



- **Differences between hinge and logistic loss**: differentiability/sparsity

# Image annotation ⇒ multi-class classification

# Losses for multi-label classification (Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004)

- **Two main strategies** for $k$ classes (with unclear winners)

1. Using existing binary classifiers (efficient code!) $+$ voting schemes
   - "one-vs-rest" : learn $k$ classifiers on the entire data
   - "one-vs-one" : learn $k(k-1)/2$ classifiers on portions of the data

# Losses for multi-label classification - Linear predictors

- Using binary classifiers (left: "one-vs-rest", right: "one-vs-one")

# Losses for multi-label classification (Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004)

- **Two main strategies** for $k$ classes (with unclear winners)

  1. Using existing binary classifiers (efficient code!) + voting schemes
     - "one-vs-rest" : learn $k$ classifiers on the entire data
     - "one-vs-one" : learn $k(k-1)/2$ classifiers on portions of the data
  2. Dedicated loss functions for prediction using $\arg\max_{i \in \{1,...,k\}} f_i(x)$
     - Softmax regression: $\text{loss} = -\log(e^{f_y(x)} / \sum_{i=1}^{k} e^{f_i(x)})$
     - Multi-class SVM - 1: $\text{loss} = \sum_{i=1}^{k} (1 + f_i(x) - f_y(x))_+$
     - Multi-class SVM - 2: $\text{loss} = \max_{i \in \{1,...,k\}} (1 + f_i(x) - f_y(x))_+$

- Different strategies do not consider same space of predictors

- Calibration of the softmax loss: $p(y|x) = \dfrac{e^{f_y(x)}}{\sum_{i=1}^{k} e^{f_i(x)}}$

# Losses for multi-label classification - Linear predictors

- Using binary classifiers (left: "one-vs-rest", right: "one-vs-one")



- Dedicated loss function

# Image retrieval ⇒ ranking

# Image retrieval ⇒ outlier/novelty detection

# Losses for ther tasks

- **Outlier detection** (Schölkopf et al., 2001; Vert and Vert, 2006)

  – one-class SVM: learn only with positive examples

- **Ranking**

  – simple trick: transform into learning on pairs (Herbrich et al., 2000), i.e., predict $\{x > y\}$ or $\{x \leqslant y\}$
  – More general "structured output methods" (Joachims, 2002)

- **General structured outputs**

  – Very active topic in machine learning and computer vision
  – see, e.g., Taskar (2005) and Christoph Lampert's course

# Dealing with asymmetric cost or unbalanced data in binary classification

- Two cases with similar issues:

  - Asymmetric cost (e.g., spam filterting, detection)
  - Unbalanced data, e.g., lots of negative examples in detection

- **One number is not enough to characterize the asymmetric properties**

  - ROC curves (Flach, 2003) – cf. precision-recall curves

- Training using asymmetric losses (Bach et al., 2006)

$$\min_{f \in \mathcal{F}} \quad C_+ \sum_{i, y_i = 1} \ell(y_i f(x_i)) + C_- \sum_{i, y_i = -1} \ell(y_i f(x_i)) + \|f\|^2$$

- Natural balance: $C_+ n_+ = C_- n_-$

# ROC curves

- ROC plane $(u, v)$

- $u =$ proportion of <span style="color:red">false positives</span> $= P(f(x) = 1 | y = -1)$

- $v =$ proportion of <span style="color:red">true positives</span> $= P(f(x) = 1 | y = 1)$

- Plot a set of classifiers $f_\gamma(x)$ for $\gamma \in \mathbb{R}$ (ex: $\gamma =$ constant term)

- Equi-cost curves and convex hulls

# ROC curves

- ROC plane $(u, v)$

- $u$ = proportion of false positives = $P(f(x) = 1 | y = -1)$

- $v$ = proportion of true positives = $P(f(x) = 1 | y = 1)$

- Plot a set of classifiers $f_\gamma(x)$ for $\gamma \in \mathbb{R}$ (ex: $\gamma$ = constant term)

- Equi-cost curves and convex hulls

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Regularizations

- Main goal: avoid overfitting (see, e.g. Hastie et al., 2001)

- Two main lines of work:

  1. Use Hilbertian (RKHS) norms
     - Non parametric supervised learning and kernel methods
     - Well developped theory (Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004; Wahba, 1990)
  2. Use "sparsity inducing" norms
     - main example: $\ell_1$-norm $\|w\|_1 = \sum_{i=1}^{p} |w_i|$
     - Perform model selection as well as regularization
     - Theory "in the making"

- **Goal of (this part of) the course: Understand how and when to use these different norms**

# Kernel methods for machine learning

- **Definition**: given a set of objects $\mathcal{X}$, a <span style="color:red">positive definite kernel</span> is a symmetric function $k(x, x')$ such that for all finite sequences of points $x_i \in \mathcal{X}$ and $\alpha_i \in \mathbb{R}$, $i = 1, \ldots, n$,

$$\sum_{i,j=1}^{n} \alpha_i \alpha_j k(x_i, x_j) \geqslant 0$$

(i.e., the matrix $(k(x_i, x_j))_{1 \leqslant i,j \leqslant n}$ is symmetric positive semi-definite)

- Main example: $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$

# Kernel methods for machine learning

- **Definition**: given a set of objects $\mathcal{X}$, a positive definite kernel is a symmetric function $k(x, x')$ such that for all finite sequences of points $x_i \in \mathcal{X}$ and $\alpha_i \in \mathbb{R}$, $i = 1, \ldots, n$,

$$\sum_{i,j=1}^{n} \alpha_i \alpha_j k(x_i, x_j) \geqslant 0$$

(i.e., the matrix $(k(x_i, x_j))_{1 \leqslant i,j \leqslant n}$ is symmetric positive semi-definite)

- **Aronszajn theorem** (Aronszajn, 1950): $k$ is a positive definite kernel if and only if there exists a Hilbert space $\mathcal{F}$ and a mapping $\Phi : \mathcal{X} \mapsto \mathcal{F}$ such that

$$\forall (x, x') \in \mathcal{X}^2, \ k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{F}}$$

- $\mathcal{X}$ = "input space", $\mathcal{F}$ = "feature space", $\Phi$ = "feature map"

- Functional view: reproducing kernel Hilbert spaces

# Classical kernels: kernels on vectors $x \in \mathbb{R}^d$

- **Linear** kernel $k(x, y) = x^\top y$

  - $\Phi(x) = x$

- **Polynomial** kernel $k(x, y) = (1 + x^\top y)^d$

  - $\Phi(x) = $ monomials
  - $(1 + x^\top y)^d = \displaystyle\sum_{\alpha_1 + \cdots + \alpha_k \leqslant d} \binom{d}{\alpha_1, \ldots, \alpha_k} (x_1 y_1)^{\alpha_1} \cdots (x_k y_k)^{\alpha_k}$

- **Gaussian** kernel $k(x, y) = \exp(-\alpha \|x - y\|^2)$

  - $\Phi(x) = $??
  - From linear classifiers ($\alpha$ small) to nearest-neighbor ($\alpha$ large)

# Reproducing kernel Hilbert spaces

- Assume $k$ is a positive definite kernel on $\mathcal{X} \times \mathcal{X}$

- **Aronszajn theorem** (1950): there exists a Hilbert space $\mathcal{F}$ and a mapping $\Phi : \mathcal{X} \mapsto \mathcal{F}$ such that

$$\forall (x, x') \in \mathcal{X}^2, \ k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$$

- $\mathcal{X} =$ "input space", $\mathcal{F} =$ "feature space", $\Phi =$ "feature map"

- RKHS: particular instantiation of $\mathcal{F}$ as a function space
  - $\Phi(x) = k(\cdot, x)$
  - function evaluation $\boxed{f(x) = \langle f, \Phi(x) \rangle}$
  - reproducing property: $k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle$

- Notations : $f(x) = \langle f, \Phi(x) \rangle = f^{\top} \Phi(x), \ \|f\|^2 = \langle f, f \rangle$

# Classical kernels: kernels on vectors $x \in \mathbb{R}^d$

- **Linear** kernel $k(x, y) = x^\top y$

  – Linear functions

- **Polynomial** kernel $k(x, y) = (1 + x^\top y)^d$

  – Polynomial functions

- **Gaussian** kernel $k(x, y) = \exp(-\alpha \|x - y\|^2)$

  – Smooth functions

# Classical kernels: kernels on vectors $x \in \mathbb{R}^d$

- **Linear** kernel $k(x, y) = x^\top y$

  - Linear functions

- **Polynomial** kernel $k(x, y) = (1 + x^\top y)^d$

  - Polynomial functions

- **Gaussian** kernel $k(x, y) = \exp(-\alpha \|x - y\|^2)$

  - Smooth functions

- **Parameter selection? Structured domain?**

  - Data are not always vectors!

# Regularization and representer theorem

- Data: $x_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$, $i = 1, \ldots, n$, kernel $k$ (with RKHS $\mathcal{F}$)

- Minimize with respect to $f$: $\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \ell(y_i, f^\top \Phi(x_i)) + \frac{\lambda}{2} \|f\|^2}$

- No assumptions on cost $\ell$ or $n$

- **Representer theorem** (Kimeldorf and Wahba, 1971): optimum is reached for weights of the form
$$f = \sum_{j=1}^{n} \alpha_j \Phi(x_j) = \sum_{j=1}^{n} \alpha_j k(\cdot, x_j)$$

- PROOF (two lines)

# Regularization and representer theorem

- Data: $x_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$, $i = 1, \ldots, n$, kernel $k$ (with RKHS $\mathcal{F}$)

- Minimize with respect to $f$: $\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \ell(y_i, f^\top \Phi(x_i)) + \frac{\lambda}{2}\|f\|^2}$

- <span style="color:red">No assumptions on cost $\ell$ or $n$</span>

- **Representer theorem** (Kimeldorf and Wahba, 1971): optimum is reached for weights of the form
$$f = \sum_{j=1}^{n} \alpha_j \Phi(x_j) = \sum_{j=1}^{n} \alpha_j k(\cdot, x_j)$$

- $\alpha \in \mathbb{R}^n$ <span style="color:red">dual parameters</span>, $K \in \mathbb{R}^{n \times n}$ <span style="color:red">kernel matrix</span>:
$$K_{ij} = \Phi(x_i)^\top \Phi(x_j) = k(x_i, x_j)$$

- Equivalent problem: $\boxed{\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^{n} \ell(y_i, (K\alpha)_i) + \frac{\lambda}{2}\alpha^\top K \alpha}$

# Kernel trick and modularity

- **Kernel trick**: any algorithm for finite-dimensional vectors that only uses pairwise dot-products can be applied in the feature space.

  - Replacing dot-products by kernel functions
  - Implicit use of (very) large feature spaces
  - Linear to non-linear learning methods

# Kernel trick and modularity

- **Kernel trick**: any algorithm for finite-dimensional vectors that only uses pairwise dot-products can be applied in the feature space.

  – Replacing dot-products by kernel functions
  – Implicit use of (very) large feature spaces
  – Linear to non-linear learning methods

- **Modularity** of kernel methods

  1. Work on new algorithms and theoretical analysis
  2. Work on new kernels for specific data types

# Representer theorem and convex duality

- The parameters $\alpha \in \mathbb{R}^n$ may also be interpreted as Lagrange multipliers

- Assumption: cost function is convex, $\varphi_i(u_i) = \ell(y_i, u_i)$

- Primal problem: $\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2}\|f\|^2}$

- What about the constant term $b$? replace $\Phi(x)$ by $(\Phi(x), c)$, $c$ large

| | $\varphi_i(u_i)$ |
|---|---|
| **LS regression** | $\frac{1}{2}(y_i - u_i)^2$ |
| **Logistic regression** | $\log(1 + \exp(-y_i u_i))$ |
| **SVM** | $(1 - y_i u_i)_+$ |

# Representer theorem and convex duality
## Proof

- **Primal** problem:
$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2} \|f\|^2$$

- Define $\varphi_i^*(v_i) = \max_{u_i \in \mathbb{R}} v_i u_i - \varphi_i(u_i)$ as the *Fenchel conjugate* of $\varphi_i$

- Main trick: introduce constraint $u_i = f^\top \Phi(x_i)$ and associated Lagrange multipliers $\alpha_i$

- Lagrangian $\mathcal{L}(\alpha, f) = \sum_{i=1}^{n} \varphi_i(u_i) + \frac{\lambda}{2}\|f\|^2 + \lambda \sum_{i=1}^{n} \alpha_i(u_i - f^\top \Phi(x_i))$

  - Maximize with respect to $u_i \Rightarrow$ term of the form $-\varphi_i^*(-\lambda \alpha_i)$
  - Maximize with respect to $f \Rightarrow f = \sum_{i=1}^{n} \alpha_i \Phi(x_i)$ and term of the form $-\frac{\lambda}{2}\|\sum_{i=1}^{n} \alpha_i \Phi(x_i)\|^2 = -\frac{\lambda}{2}\alpha^\top K \alpha$

# Representer theorem and convex duality

- Assumption: cost function is <span style="color:red">convex</span> $\varphi_i(u_i) = \ell(y_i, u_i)$

- <span style="color:red">Primal</span> problem: 
$$\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^n \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2}\|f\|^2}$$

- <span style="color:red">Dual</span> problem: 
$$\boxed{\max_{\alpha \in \mathbb{R}^n} -\sum_{i=1}^n \varphi_i^*(-\lambda \alpha_i) - \frac{\lambda}{2}\alpha^\top K \alpha}$$

  where $\varphi_i^*(v_i) = \max_{u_i \in \mathbb{R}} v_i u_i - \varphi_i(u_i)$ is the Fenchel conjugate of $\varphi_i$

- Strong duality

- Relationship between primal and dual variables (at optimum):
$$f = \sum_{i=1}^n \alpha_i \Phi(x_i)$$

- NB: adding constant term $b \Leftrightarrow$ add constraints $\sum_{i=1}^n \alpha_i = 0$

# Supervised kernel methods (2-norm regularization)

Primal problem $\quad\min_{f \in \mathcal{F}} \left( \sum_i \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2}\|f\|^2 \right)$

Dual problem $\quad\max_{\alpha \in \mathbb{R}^n} \left( -\sum_i \varphi_i^*(\lambda \alpha_i) - \frac{\lambda}{2}\alpha^\top K \alpha \right)$

Optimality conditions $\quad f = \sum_{i=1}^{n} \alpha_i \Phi(x_i)$

- Assumptions on loss $\varphi_i$:

  - $\varphi_i(u)$ convex
  - $\varphi_i^*$ Fenchel conjugate of $\varphi_i$, i.e., $\varphi_i^*(v) = \max_{u \in \mathbb{R}}(vu - \varphi_i(u))$

|  | $\varphi_i(u_i)$ | $\varphi_i^*(v)$ |
|---|---|---|
| **LS regression** | $\frac{1}{2}(y_i - u_i)^2$ | $\frac{1}{2}v^2 + vy_i$ |
| **Logistic regression** | $\log(1 + \exp(-y_i u_i))$ | $(1+vy_i)\log(1+vy_i)$ $-vy_i \log(-vy_i)$ |
| **SVM** | $(1 - y_i u_i)_+$ | $vy_i \times 1_{-vy_i \in [0,1]}$ |

# Particular case of the support vector machine

- Primal problem: $\boxed{\min\limits_{f\in\mathcal{F}} \sum_{i=1}^{n}(1 - y_i f^\top \Phi(x_i))_+ + \frac{\lambda}{2}\|f\|^2}$

- Dual problem: $\boxed{\max\limits_{\alpha\in\mathbb{R}^n}\left(-\sum_i \lambda\alpha_i y_i \times 1_{-\lambda\alpha_i y_i\in[0,1]} - \frac{\lambda}{2}\alpha^\top K\alpha\right)}$

- Dual problem (by change of variable $\alpha \leftarrow -\operatorname{Diag}(y)\alpha$ and $C = 1/\lambda$):

$$\boxed{\max\limits_{\alpha\in\mathbb{R}^n,\ 0\leqslant\alpha\leqslant C} \sum_{i=1}^{n}\alpha_i - \frac{1}{2}\alpha^\top \operatorname{Diag}(y)K\operatorname{Diag}(y)\alpha}$$

# Particular case of the support vector machine

- **Primal** problem: $\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (1 - y_i f^\top \Phi(x_i))_+ + \frac{\lambda}{2} \|f\|^2}$

- **Dual** problem:

$$\boxed{\max_{\alpha \in \mathbb{R}^n, \ 0 \leqslant \alpha \leqslant C} \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\alpha^\top \operatorname{Diag}(y) K \operatorname{Diag}(y)\alpha}$$

# Particular case of the support vector machine

- Primal problem: $\boxed{\min_{f \in \mathcal{F}} \sum_{i=1}^{n} (1 - y_i f^\top \Phi(x_i))_+ + \frac{\lambda}{2} \|f\|^2}$

- Dual problem:

$$\boxed{\max_{\alpha \in \mathbb{R}^n,\ 0 \leqslant \alpha \leqslant C} \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\alpha^\top \operatorname{Diag}(y) K \operatorname{Diag}(y)\alpha}$$

- What about the traditional picture?

# Losses for classification (Shawe-Taylor and Cristianini, 2004)

- **Label** : $y \in \{-1, 1\}$, prediction $\hat{y} = \text{sign}(f(x))$

  - loss of the form $\ell(y, f(x)) = \ell(yf(x))$
  - "True" cost: $\ell(yf(x)) = 1_{yf(x)<0}$
  - Usual <span style="color:red">convex</span> costs:



- **Differences between hinge and logistic loss**: differentiability/sparsity

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Kernel ridge regression (a.k.a. spline smoothing) - I

- Data $x_1, \ldots, x_n \in \mathcal{X}$, $y_1, \ldots, y_n \in \mathbb{R}$, positive definite kernel $k$

- Least-squares

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

- **View 1**: representer theorem $\Rightarrow f = \sum_{i=1}^{n} \alpha_i k(\cdot, x_i)$

  – equivalent to

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^{n} (y_i - (K\alpha)_i)^2 + \lambda \alpha^\top K \alpha$$

  – Solution equal to $\alpha = (K + n\lambda I)^{-1} y + \varepsilon$ with $K\varepsilon = 0$
  – Unique solution $f$

# Kernel ridge regression (a.k.a. spline smoothing) - II

- Links with spline smoothing (Wahba, 1990)

- **View 2**: primal problem $\mathcal{F} \subset \mathbb{R}^d$, $\Phi \in \mathbb{R}^{n \times d}$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \|y - \Phi w\|^2 + \lambda \|w\|^2$$

- Solution equal to $w = (\Phi^\top \Phi + n\lambda I)^{-1} \Phi^\top y$

- Note that $w = \Phi^\top (\Phi \Phi^\top + n\lambda I)^{-1} y = \Phi^\top (K + n\lambda I)^{-1} y = \Phi^\top \alpha$

  – Using matrix inversion lemma

- $\Phi w$ equal to $K\alpha$

# Kernel ridge regression (a.k.a. spline smoothing) - III

- **View 3**: dual problem

    – dual problem: $\max_{\alpha \in \mathbb{R}^n} -\frac{n\lambda}{2}\|\alpha\|^2 - \alpha^\top y - \frac{1}{2}\alpha^\top K\alpha$
    – solution: $\alpha = (K + \lambda I)^{-1} y$

- Warning: same solution obtained from different point of views

    – Primal problem: linear system of size $d$
    – Dual problem: linear system of size $n$

# Losses for classification

- Usual convex costs:



- **Differences between hinge and logistic loss**: differentiability/sparsity

# Support vector machine or logistic regression?

- Predictive performance is similar

  – Logistic cost better calibrated and often easier to optimize

- Only true difference is numerical

  – SVM: sparsity in $\alpha$
  – Logistic: differentiable loss function

- Which one to use?

  – Linear kernel $\Rightarrow$ Logistic + Newton/Gradient descent
  – Linear kernel - Large scale $\Rightarrow$ Stochastic gradient descent
  – Nonlinear kernel $\Rightarrow$ SVM + dual methods or simpleSVM

# Algorithms for supervised kernel methods

- Four formulations

  1. Dual: $\max_{\alpha \in \mathbb{R}^n} - \sum_i \varphi_i^*(\lambda \alpha_i) - \frac{\lambda}{2} \alpha^\top K \alpha$
  2. Primal: $\min_{f \in \mathcal{F}} \sum_i \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2} \|f\|^2$
  3. Primal + Representer: $\min_{\alpha \in \mathbb{R}^n} \sum_i \varphi_i((K\alpha)_i) + \frac{\lambda}{2} \alpha^\top K \alpha$
  4. Convex programming

- **Best strategy depends on loss (differentiable or not) and kernel (linear or not)**

# Dual methods

- Dual problem: $\max_{\alpha \in \mathbb{R}^n} - \sum_i \varphi_i^*(\lambda \alpha_i) - \frac{\lambda}{2} \alpha^\top K \alpha$

- Main method: coordinate descent (a.k.a. sequential minimal optimization - SMO) (Platt, 1998; Bottou and Lin, 2007; Joachims, 1998)

  - Efficient when loss is piecewise quadratic (i.e., hinge $=$ SVM)
  - Sparsity may be used in the case of the SVM

- Computational complexity: between quadratic and cubic in $n$

- **Works for all kernels**

# Primal methods

- Primal problem: $\min_{f \in \mathcal{F}} \sum_i \varphi_i(f^\top \Phi(x_i)) + \frac{\lambda}{2} ||f||^2$

- Only works directly if $\Phi(x)$ may be built explicitly and has small dimension

  – Example: linear kernel in small dimensions

- Differentiable loss: gradient descent or Newton's method are very efficient in small dimensions

- **Do not use linear SVMs in small dimensions!**

- Larger scale

  – stochastic gradient descent (Shalev-Shwartz et al., 2007; Bottou and Bousquet, 2008)
  – See Zaid Harchaoui's course

# Primal methods with representer theorems

- Primal problem in $\alpha$: $\min_{\alpha \in \mathbb{R}^n} \sum_i \varphi_i((K\alpha)_i) + \frac{\lambda}{2}\alpha^\top K\alpha$

- Direct optimization in $\alpha$ poorly conditioned ($K$ has low-rank) unless Newton method is used (Chapelle, 2007)

- General kernels: use incomplete Cholesky decomposition (Fine and Scheinberg, 2001; Bach and Jordan, 2002) to obtain a square root $K = GG^\top$



$G$ of size $n \times m$, where $m \ll n$

  – "Empirical input space" of size $m$ obtained using rows of $G$
  – Running time to compute $G$: $O(m^2 n)$

# Direct convex programming

- **Convex programming toolboxes $\Rightarrow$ very inefficient!**

- May use special structure of the problem

  – e.g., SVM and sparsity in $\alpha$

- Active set method for the SVM: **SimpleSVM** (Vishwanathan et al., 2003; Loosli et al., 2005)

  – Cubic complexity in the number of support vectors

- Full regularization path for the SVM (Hastie et al., 2005; Bach et al., 2006)

  – Cubic complexity in the number of support vectors
  – May be extended to other settings (Rosset and Zhu, 2007)

# Code

- SVM and other supervised learning techniques

  `www.shogun-toolbox.org`

  `http://gaelle.loosli.fr/research/tools/simplesvm.html`

  `http://www.kyb.tuebingen.mpg.de/bs/people/spider/main.html`

  `http://ttic.uchicago.edu/~shai/code/index.html`

- $\ell^1$-penalization:

  - SPAMS (SPArse Modeling Software)

    `http://www.di.ens.fr/willow/SPAMS/`

- Multiple kernel learning:

  `asi.insa-rouen.fr/enseignants/~arakotom/code/mklindex.html`

  `www.stat.berkeley.edu/~gobo/SKMsmo.tar`

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Kernel methods - I

- Distances in the "feature space"

$$d_k(x, y)^2 = \|\Phi(x) - \Phi(y)\|_{\mathcal{F}}^2 = k(x, x) + k(y, y) - 2k(x, y)$$

- Nearest-neighbor classification/regression

  – Gaussian kernel: $k(x, x) = k(y, y) = 1 \Rightarrow$ same as in input space

# Kernel methods - II
## Simple discrimination algorithm

- Data $x_1, \ldots, x_n \in \mathcal{X}$, classes $y_1, \ldots, y_n \in \{-1, 1\}$

- Compare distances to mean of each class

- Equivalent to classifying $x$ using the sign of

$$\frac{1}{\#\{i, y_i = 1\}} \sum_{i, y_i = 1} k(x, x_i) - \frac{1}{\#\{i, y_i = -1\}} \sum_{i, y_i = -1} k(x, x_i)$$

- Proof...

- Geometric interpretation of Parzen windows

  - NB: onyl when $k$ is positive definite and pointwise positive (i.e., $\forall x, y, \ k(x, y) \geqslant 0$)

# Kernel methods - III
# Data centering

- $n$ points $x_1, \ldots, x_n \in \mathcal{X}$

- kernel matrix $K \in \mathbb{R}^n$, $K_{ij} = k(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$

- Kernel matrix of centered data $\tilde{K}_{ij} = \langle \Phi(x_i) - \mu, \Phi(x_j) - \mu \rangle$ where $\mu = \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i)$

- Formula: $\tilde{K} = \Pi_n K \Pi_n$ with $\Pi_n = I_n - \frac{E}{n}$, and $E$ constant matrix equal to 1.

- Proof...

- NB: $\mu$ is not of the form $\Phi(z)$, $z \in \mathcal{X}$ (cf. preimage problem)

# Kernel PCA

- Linear principal component analysis

  - data $x_1, \ldots, x_n \in \mathbb{R}^p$,

$$\max_{w \in \mathbb{R}^p} \frac{w^\top \hat{\Sigma} w}{w^\top w} = \max_{w \in \mathbb{R}^p} \frac{\mathrm{var}(w^\top X)}{w^\top w}$$

  - $w$ is largest eigenvector of $\hat{\Sigma}$
  - Denoising, data representation

- Kernel PCA: data $x_1, \ldots, x_n \in \mathcal{X}$, p.d. kernel $k$

  - View 1: $\max_{w \in \mathcal{F}} \dfrac{\mathrm{var}(\langle \Phi(X), w \rangle)}{w^\top w}$ \qquad View 2: $\max_{f \in \mathcal{F}} \dfrac{\mathrm{var}(f(X))}{\|f\|_{\mathcal{F}}^2}$
  - Solution: $f, w = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$ and $\alpha$ first eigenvector of $\tilde{K} = \Pi_n K \Pi_n$
  - Interpretation in terms of covariance operators

# Denoising with kernel PCA (From Schölkopf, 2005)

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Kernel design

- Principle: **kernel on $\mathcal{X}$ = space of functions on $\mathcal{X}$ + norm**

- Two main design principles

  1. Constructing kernels from kernels by algebraic operations
  2. Using usual algebraic/numerical tricks to perform efficient kernel computations with very high-dimensional feature spaces

- Operations: $k_1(x, y) = \langle \Phi_1(x), \Phi_1(y) \rangle$, $k_2(x, y) = \langle \Phi_2(x), \Phi_2(y) \rangle$

  - **Sum = concatenation of feature spaces**:

$$k_1(x, y) + k_2(x, y) = \left\langle \begin{pmatrix} \Phi_1(x) \\ \Phi_2(x) \end{pmatrix}, \begin{pmatrix} \Phi_1(y) \\ \Phi_2(y) \end{pmatrix} \right\rangle$$

  - **Product = tensor product of feature spaces**:

$$k_1(x, y) k_2(x, y) = \left\langle \Phi_1(x) \Phi_2(x)^\top, \Phi_1(y) \Phi_2(y)^\top \right\rangle$$

# Classical kernels: kernels on vectors $x \in \mathbb{R}^d$

- **Linear** kernel $k(x, y) = x^\top y$

  – Linear functions

- **Polynomial** kernel $k(x, y) = (1 + x^\top y)^d$

  – Polynomial functions

- **Gaussian** kernel $k(x, y) = \exp(-\alpha \|x - y\|^2)$

  – Smooth functions

- Data are not always vectors!

# Efficient ways of computing large sums

- Goal: $\Phi(x) \in \mathbb{R}^p$ high-dimensional, compute $\displaystyle\sum_{i=1}^{p} \Phi_i(x)\Phi_i(y)$ **in** $o(p)$

- **Sparsity**: many $\Phi_i(x)$ equal to zero (example: pyramid match kernel)

- **Factorization and recursivity**: replace sums of many products by product of few sums (example: polynomial kernel, graph kernel)

$$(1 + x^\top y)^d = \sum_{\alpha_1 + \cdots + \alpha_k \leqslant d} \binom{d}{\alpha_1, \ldots, \alpha_k} (x_1 y_1)^{\alpha_1} \cdots (x_k y_k)^{\alpha_k}$$

# Kernels over (labelled) sets of points

- Common situation in computer vision (e.g., interest points)

- Simple approach: compute averages/histograms of certain features

  - valid kernels over histograms $h$ and $h'$ (Hein and Bousquet, 2004)
  - **intersection**: $\sum_i \min(h_i, h'_i)$, **chi-square**: $\exp\left(-\alpha \sum_i \frac{(h_i - h'_i)^2}{h_i + h'_i}\right)$

# Kernels over (labelled) sets of points

- Common situation in computer vision (e.g., interest points)

- Simple approach: compute averages/histograms of certain features

  - valid kernels over histograms $h$ and $h'$ (Hein and Bousquet, 2004)
  - **intersection**: $\sum_i \min(h_i, h'_i)$, **chi-square**: $\exp\left(-\alpha \sum_i \frac{(h_i - h'_i)^2}{h_i + h'_i}\right)$

- Pyramid match (Grauman and Darrell, 2007): efficiently introducing localization

  - Form a regular pyramid on top of the image
  - Count the number of common elements in each bin
  - Give a weight to each bin
  - Many bins but most of them are empty
    $\Rightarrow$ use sparsity to compute kernel efficiently

# Pyramid match kernel
## (Grauman and Darrell, 2007; Lazebnik et al., 2006)

- Two sets of points



- Counting matches at several scales: 7, 5, 4

# Kernels from segmentation graphs

- Goal of segmentation: extract objects of interest

- Many methods available, ....

  - ... but, rarely find the object of interest entirely

- Segmentation graphs

  - Allows to work on "more reliable" over-segmentation
  - Going to a <span style="color:red">large square grid (millions of pixels)</span> to a <span style="color:red">small graph (dozens or hundreds of regions)</span>

- How to build a kernel over segmenation graphs?

  - NB: more generally, kernelizing existing representations?

# Segmentation by watershed transform (Meyer, 2001)



image        gradient        watershed

287 segments        64 segments        10 segments

# Segmentation by watershed transform (Meyer, 2001)

image

gradient

watershed

287 segments

64 segments

10 segments

# Image as a segmentation graph

- Labelled undirected graph

  - Vertices: connected segmented regions
  - Edges: between spatially neighboring regions
  - Labels: region pixels

# Image as a segmentation graph

- <span style="color:red">Labelled undirected graph</span>

  - <span style="color:blue">Vertices</span>: connected segmented regions
  - <span style="color:blue">Edges</span>: between spatially neighboring regions
  - <span style="color:blue">Labels</span>: region pixels

- Difficulties

  - Extremely high-dimensional labels
  - Planar undirected graph
  - Inexact matching

- **<span style="color:red">Graph kernels</span>** (Gärtner et al., 2003; Kashima et al., 2004; Harchaoui and Bach, 2007) provide an elegant and efficient solution

# Kernels between structured objects
## Strings, graphs, etc… (Shawe-Taylor and Cristianini, 2004)

- Numerous applications (text, bio-informatics, speech, vision)

- Common design principle: **enumeration of subparts** (Haussler, 1999; Watkins, 1999)

  - Efficient for strings
  - Possibility of gaps, partial matches, very efficient algorithms

- Most approaches fail for general graphs (even for undirected trees!)

  - NP-Hardness results (Ramon and Gärtner, 2003)
  - Need specific set of subparts

# Paths and walks

- Given a graph $G$,

  - A path is a sequence of distinct neighboring vertices
  - A walk is a sequence of neighboring vertices

- Apparently similar notions

# Walks

# Walk kernel (Kashima et al., 2004; Borgwardt et al., 2005)

- $\mathcal{W}_{\mathbf{G}}^p$ (resp. $\mathcal{W}_{\mathbf{H}}^p$) denotes the set of walks of length $p$ in $\mathbf{G}$ (resp. $\mathbf{H}$)

- Given *basis kernel* on labels $k(\ell, \ell')$

- *$p$-th order walk kernel*:

$$k_{\mathcal{W}}^p(\mathbf{G}, \mathbf{H}) = \sum_{\substack{(r_1, \ldots, r_p) \in \mathcal{W}_{\mathbf{G}}^p \\ (s_1, \ldots, s_p) \in \mathcal{W}_{\mathbf{H}}^p}} \prod_{i=1}^p k(\ell_{\mathbf{G}}(r_i), \ell_{\mathbf{H}}(s_i)).$$

# Dynamic programming for the walk kernel (Harchaoui and Bach, 2007)

- Dynamic programming in $O(p d_{\mathbf{G}} d_{\mathbf{H}} n_{\mathbf{G}} n_{\mathbf{H}})$

- $k_{\mathcal{W}}^p(\mathbf{G}, \mathbf{H}, r, s) = $ sum restricted to walks starting at $r$ and $s$

- recursion between $p-1$-th walk and $p$-th walk kernel

$$k_{\mathcal{W}}^p(\mathbf{G}, \mathbf{H}, r, s) = k(\ell_{\mathbf{G}}(r), \ell_{\mathbf{H}}(s)) \sum_{\substack{r' \in \mathcal{N}_{\mathbf{G}}(r) \\ s' \in \mathcal{N}_{\mathbf{H}}(s)}} k_{\mathcal{W}}^{p-1}(\mathbf{G}, \mathbf{H}, r', s').$$

# Dynamic programming for the walk kernel (Harchaoui and Bach, 2007)

- Dynamic programming in $O(pd_{\mathbf{G}}d_{\mathbf{H}}n_{\mathbf{G}}n_{\mathbf{H}})$

- $k^p_{\mathcal{W}}(\mathbf{G}, \mathbf{H}, r, s) = $ sum restricted to walks starting at $r$ and $s$

- recursion between $p-1$-th walk and $p$-th walk kernel

$$k^p_{\mathcal{W}}(\mathbf{G}, \mathbf{H}, r, s) = k(\ell_{\mathbf{G}}(r), \ell_{\mathbf{H}}(s)) \sum_{\substack{r' \in \mathcal{N}_{\mathbf{G}}(r) \\ s' \in \mathcal{N}_{\mathbf{H}}(s)}} k^{p-1}_{\mathcal{W}}(\mathbf{G}, \mathbf{H}, r', s')$$

- Kernel obtained as $k^{p,\alpha}_{\mathcal{T}}(\mathbf{G}, \mathbf{H}) = \sum_{r \in \mathcal{V}_{\mathbf{G}}, s \in \mathcal{V}_{\mathbf{H}}} k^{p,\alpha}_{\mathcal{T}}(\mathbf{G}, \mathbf{H}, r, s)$

# Extensions of graph kernels

- Main principle: **compare all possible subparts of the graphs**

- Going from paths to subtrees

  - Extension of the concept of walks $\Rightarrow$ tree-walks (Ramon and Gärtner, 2003)

- Similar dynamic programming recursions (Harchaoui and Bach, 2007)

- Need to play around with subparts to obtain efficient recursions

  - NB: Do we actually need positive definiteness?

# Performance on Corel14 (Harchaoui and Bach, 2007)

- Corel14: 1400 natural images with 14 classes

# Performance on Corel14 (Harchaoui & Bach, 2007) Error rates

- Histogram kernels (**H**)

- Walk kernels (**W**)

- Tree-walk kernels (**TW**)

- Weighted tree-walks (**wTW**)

- MKL (**M**)



Performance comparison on Corel14

- Hyperparameter selection using cross-validation

# Kernel methods - Summary

- **Kernels and representer theorems**

  – Clear distinction between representation/algorithms

- **Algorithms**

  – Two formulations (primal/dual)
  – Logistic or SVM?

- **Kernel design**

  – Very large feature spaces with efficient kernel evaluations

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Supervised learning and regularization

- Data: $x_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$, $i = 1, \ldots, n$

- Minimize with respect to function $f : \mathcal{X} \to \mathcal{Y}$:

$$\sum_{i=1}^{n} \ell(y_i, f(x_i)) \qquad + \qquad \frac{\lambda}{2}\|f\|^2$$

$$\text{Error on data} \qquad + \qquad \text{Regularization}$$

$$\text{Loss \& function space ?} \qquad \qquad \text{Norm ?}$$

- Two theoretical/algorithmic issues:

1. Loss
2. **Function space / norm**

# Regularizations

- **Main goal: avoid overfitting**

- **Two main lines of work**:

  1. Euclidean and Hilbertian norms (i.e., $\ell_2$-norms)
     - Possibility of non linear predictors
     - Non parametric supervised learning and kernel methods
     - Well developped theory and algorithms (see, e.g., Wahba, 1990; Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004)

# Regularizations

- **Main goal: avoid overfitting**

- **Two main lines of work**:

  1. Euclidean and Hilbertian norms (i.e., $\ell_2$-norms)
     - Possibility of non linear predictors
     - Non parametric supervised learning and kernel methods
     - Well developped theory and algorithms (see, e.g., Wahba, 1990; Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004)
  2. Sparsity-inducing norms
     - Usually restricted to linear predictors on vectors $f(x) = w^\top x$
     - Main example: $\ell_1$-norm $\|w\|_1 = \sum_{i=1}^p |w_i|$
     - Perform model selection as well as regularization
     - **Theory and algorithms "in the making"**

# $\ell_2$-**norm vs.** $\ell_1$-**norm**

- $\ell_1$-norms lead to interpretable models

- $\ell_2$-norms can be run implicitly with very large feature spaces

- **Algorithms**:

  - Smooth convex optimization vs. nonsmooth convex optimization

- **Theory**:

  - better predictive performance?

# $\ell_2$ vs. $\ell_1$ - Gaussian hare vs. Laplacian tortoise



- First-order methods (Fu, 1998; Wu and Lange, 2008)
- Homotopy methods (Markowitz, 1956; Efron et al., 2004)

# Why $\ell_1$-norms lead to sparsity?

- **Example 1**: quadratic problem in 1D, i.e. $\boxed{\min_{x \in \mathbb{R}} \dfrac{1}{2}x^2 - xy + \lambda|x|}$

- Piecewise quadratic function with a kink at zero

  - Derivative at $0+$: $g_+ = \lambda - y$ and $0-$: $g_- = -\lambda - y$



  - $x = 0$ is the solution iff $g_+ \geqslant 0$ and $g_- \leqslant 0$ (i.e., $|y| \leqslant \lambda$)
  - $x \geqslant 0$ is the solution iff $g_+ \leqslant 0$ (i.e., $y \geqslant \lambda$) $\Rightarrow x^* = y - \lambda$
  - $x \leqslant 0$ is the solution iff $g_- \leqslant 0$ (i.e., $y \leqslant -\lambda$) $\Rightarrow x^* = y + \lambda$

- Solution $\boxed{x^* = \mathrm{sign}(y)(|y| - \lambda)_+}$ = soft thresholding

# Why $\ell_1$-norms lead to sparsity?

- **Example 1**: quadratic problem in 1D, i.e. $$\min_{x \in \mathbb{R}} \frac{1}{2} x^2 - xy + \lambda |x|$$

- Piecewise quadratic function with a kink at zero

- Solution $\boxed{x^* = \text{sign}(y)(|y| - \lambda)_+}$ = soft thresholding

# Why $\ell_1$-norms lead to sparsity?

- **Example 2**: minimize quadratic function $Q(w)$ subject to $\|w\|_1 \leqslant T$.
  - coupled soft thresholding

- Geometric interpretation
  - NB : penalizing is "equivalent" to constraining

# $\ell_1$-norm regularization (linear setting)

- Data: covariates $x_i \in \mathbb{R}^p$, responses $y_i \in \mathcal{Y}$, $i = 1, \ldots, n$

- Minimize with respect to loadings/weights $w \in \mathbb{R}^p$:

$$J(w) = \sum_{i=1}^{n} \ell(y_i, w^\top x_i) \quad + \quad \lambda \|w\|_1$$

$$\text{Error on data} \quad + \quad \text{Regularization}$$

- Including a constant term $b$? Penalizing or constraining?

- square loss $\Rightarrow$ basis pursuit in signal processing (Chen et al., 2001), Lasso in statistics/machine learning (Tibshirani, 1996)

# First order methods for convex optimization on $\mathbb{R}^p$
## Smooth optimization

- **Gradient descent**: $w_{t+1} = w_t - \alpha_t \nabla J(w_t)$

  - with line search: search for a decent (not necessarily best) $\alpha_t$
  - fixed diminishing step size, e.g., $\alpha_t = t^{-1/2}$

- Convergence of $J(w_t)$ to $J^* = \min_{w \in \mathbb{R}^p} J(w)$ (Nesterov, 2003)

  - $f$ convex and $M$-Lipschitz: $\qquad\qquad J(w_t) - J^* = O\big(M/\sqrt{t}\big)$
  - and, differentiable with $L$-Lipschitz gradient: $J(w_t) - J^* = O\big(L/t\big)$
  - and, $J$ $\mu$-strongly convex: $\qquad\qquad J(w_t) - J^* = O\big(L \exp(-4t\frac{\mu}{L})\big)$

- $\frac{\mu}{L} =$ condition number of the optimization problem

- Coordinate descent: similar properties

- NB: "optimal scheme" $J(w_t) - J^* = O\big(L \min\{\exp(-4t\sqrt{\mu/L}), t^{-2}\}\big)$

# First-order methods for convex optimization on $\mathbb{R}^p$
## Non smooth optimization

- First-order methods for non differentiable objective

  - Subgradient descent: $w_{t+1} = w_t - \alpha_t g_t$, with $g_t \in \partial J(w_t)$
    - with exact line search: not always convergent (see counter-example)
    - diminishing step size, e.g., $\alpha_t = a(t+b)^{-1/2}$: convergent
  - Coordinate descent: not always convergent (show counter-example)

- Convergence rates ($J$ convex and $M$-Lipschitz): $J(w_t) - J^* = O\left(\frac{M}{\sqrt{t}}\right)$

# Counter-example
## Coordinate descent for nonsmooth objectives

# Counter-example (Bertsekas, 1995)
# Steepest descent for nonsmooth objectives

- $q(x_1, x_2) = \begin{cases} -5(9x_1^2 + 16x_2^2)^{1/2} & \text{if } x_1 > |x_2| \\ -(9x_1 + 16|x_2|)^{1/2} & \text{if } x_1 \leqslant |x_2| \end{cases}$

- Steepest descent starting from any $x$ such that $x_1 > |x_2| > (9/16)^2|x_1|$

# Second order methods

- Differentiable case

  - Newton: $w_{t+1} = w_t - \alpha_t H_t^{-1} g_t$
    - Traditional: $\alpha_t = 1$, but non globally convergent
    - globally convergent with line search for $\alpha_t$ (see Boyd, 2003)
    - $O(\log \log(1/\varepsilon))$ (slower) iterations
  - Quasi-newton methods (see Bonnans et al., 2003)

- Non differentiable case (interior point methods)

  - Smoothing of problem + second order methods
    * See example later and (Boyd, 2003)
    * Theoretically $O(\sqrt{p})$ Newton steps, usually $O(1)$ Newton steps

# Regularized problems - Proximal methods

- Gradient descent as a proximal method (differentiable functions)

  - $w_{t+1} = \arg\min\limits_{w \in \mathbb{R}^p} J(w_t) + (w - w_t)^\top \nabla J(w_t) + \dfrac{L}{2}\|w - w_t\|_2^2$
  - $w_{t+1} = w_t - \dfrac{1}{L}\nabla J(w_t)$

- Problems of the form: $\boxed{\min\limits_{w \in \mathbb{R}^p} J(w) + \lambda \Omega(w)}$

  - $w_{t+1} = \arg\min\limits_{w \in \mathbb{R}^p} J(w_t) + (w - w_t)^\top \nabla J(w_t) + \lambda \Omega(w) + \dfrac{L}{2}\|w - w_t\|_2^2$
  - Thresholded gradient descent

- Similar convergence rates than smooth optimization

  - Acceleration methods (Nesterov, 2007; Beck and Teboulle, 2009)
  - **depends on the condition number of the loss**

# Piecewise linear paths

# Algorithms for $\ell_1$-norms (square loss): Gaussian hare vs. Laplacian tortoise



- Coordinate descent: $O(pn)$ per iterations for $\ell_1$ and $\ell_2$

- "Exact" algorithms: $O(kpn)$ for $\ell_1$ **vs.** $O(p^2n)$ for $\ell_2$

- See Bach, Jenatton, Mairal, and Obozinski (2011)

# Additional methods - Softwares

- Many contributions in signal processing, optimization, machine learning

  - Extensions to stochastic setting (Bottou and Bousquet, 2008)

- Extensions to other sparsity-inducing norms

  - Computing proximal operator
  - See small `www.di.ens.fr/~fbach/bach_jenatton_mairal_obozinski_FOT.pdf`

- **Softwares**

  - Many available codes
  - SPAMS (SPArse Modeling Software) `http://www.di.ens.fr/willow/SPAMS/`

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Lasso - Two main recent theoretical results

1. **Support recovery condition** (Zhao and Yu, 2006; Wainwright, 2006; Zou, 2006; Yuan and Lin, 2007): the Lasso is sign-consistent if and only if

$$\|\mathbf{Q_{J^cJ}}\mathbf{Q_{JJ}^{-1}}\mathrm{sign}(\mathbf{w_J})\|_\infty \leqslant 1,$$

where $\mathbf{Q} = \lim_{n\to+\infty} \frac{1}{n}\sum_{i=1}^n x_i x_i^\top \in \mathbb{R}^{p\times p}$ and $\mathbf{J} = \mathrm{Supp}(\mathbf{w})$

# Lasso - Two main recent theoretical results

1. **Support recovery condition** (Zhao and Yu, 2006; Wainwright, 2006; Zou, 2006; Yuan and Lin, 2007): the Lasso is sign-consistent if and only if

$$\|\mathbf{Q}_{\mathbf{J}^c\mathbf{J}}\mathbf{Q}_{\mathbf{J}\mathbf{J}}^{-1}\mathrm{sign}(\mathbf{w_J})\|_{\infty} \leqslant 1,$$

where $\mathbf{Q} = \lim_{n \to +\infty} \frac{1}{n}\sum_{i=1}^{n} x_i x_i^\top \in \mathbb{R}^{p \times p}$ and $\mathbf{J} = \mathrm{Supp}(\mathbf{w})$

2. **Exponentially many irrelevant variables** (Zhao and Yu, 2006; Wainwright, 2006; Bickel et al., 2009; Lounici, 2008; Meinshausen and Yu, 2008): under appropriate assumptions, consistency is possible as long as

$$\log p = O(n)$$

# Alternative sparse methods
## Greedy methods

- Forward selection

- Forward-backward selection

- Non-convex method

  - Harder to analyze
  - Simpler to implement
  - Problems of stability

- Positive theoretical results (Zhang, 2009, 2008)

  - Similar sufficient conditions than for the Lasso

# Comparing Lasso and other strategies for linear regression

- Compared methods to reach the least-square solution

  - Ridge regression: $\min\limits_{w\in\mathbb{R}^p} \dfrac{1}{2}\|y - Xw\|_2^2 + \dfrac{\lambda}{2}\|w\|_2^2$

  - Lasso: $\min\limits_{w\in\mathbb{R}^p} \dfrac{1}{2}\|y - Xw\|_2^2 + \lambda\|w\|_1$

  - Forward greedy:
    * Initialization with empty set
    * Sequentially add the variable that best reduces the square loss

- Each method builds a path of solutions from 0 to ordinary least-squares solution

- Regularization parameters selected on the test set

# Simulation results

- i.i.d. Gaussian design matrix, $k = 4$, $n = 64$, $p \in [2, 256]$, SNR $= 1$
- Note stability to non-sparsity and variability



Sparse

# Simulation results

- i.i.d. Gaussian design matrix, $k = 4$, $n = 64$, $p \in [2, 256]$, SNR $= 1$
- Note stability to non-sparsity and variability



Sparse

Rotated (non sparse)

# **Summary**
# $\ell_1$-**norm regularization**

- $\ell_1$-norm regularization leads to **nonsmooth optimization problems**

  - analysis through directional derivatives or subgradients
  - optimization may or may not take advantage of sparsity

- $\ell_1$-norm regularization allows **high-dimensional inference**

- Interesting problems for $\ell_1$-regularization

  - Stable variable selection
  - Weaker sufficient conditions (for weaker results)
  - Estimation of regularization parameter (all bounds depend on the unknown noise variance $\sigma^2$)

# Extensions

- **Sparse methods are not limited to the square loss**

  – logistic loss: algorithms (Beck and Teboulle, 2009) and theory (Van De Geer, 2008; Bach, 2009)

- **Sparse methods are not limited to supervised learning**

  – Learning the structure of Gaussian graphical models (Meinshausen and Bühlmann, 2006; Banerjee et al., 2008)
  – Sparsity on matrices (last part of the course)
  – See Jean Ponce's course

- **Sparse methods are not limited to variable selection in a linear model**

  – **See next part of the course**

# Course outline

1. **Losses for particular machine learning tasks**

   • Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   • Kernels and representer theorem
   • Convex duality, optimization and algorithms
   • Kernel methods
   • Kernel design

3. **Regularization by sparsity-inducing norms**

   • $\ell_1$-norm regularization
   • Theoretical results
   • Multiple kernel learning
   • Learning on matrices

# Penalization with grouped variables (Yuan and Lin, 2006)

- Assume that $\{1, \ldots, p\}$ is **partitioned** into $m$ groups $G_1, \ldots, G_m$

- Penalization by $\sum_{i=1}^{m} \|w_{G_i}\|_2$, often called $\ell_1$-$\ell_2$ norm

- Induces group sparsity

  – Some groups entirely set to zero
  – no zeros within groups

$G_1$

$G_2$

$G_3$

# Unit norm balls
## Geometric interpretation

$$\|w\|_2 \qquad\qquad \|w\|_1 \qquad\qquad \sqrt{w_1^2 + w_2^2} + |w_3|$$

# Linear vs. non-linear methods

- All methods in this course are **linear in the parameters**

- By replacing $x$ by features $\Phi(x)$, they can be made **non linear in the data**

- **Implicit vs. explicit features**

  - $\ell_1$-norm: explicit features
  - $\ell_2$-norm: representer theorem allows to consider implicit features if their dot products can be computed easily (kernel methods)

# **Multiple kernel learning (MKL)**
# **(Lanckriet et al., 2004b; Bach et al., 2004a)**

- Sparse methods are linear!

- Sparsity with non-linearities

  - replace $f(x) = \sum_{j=1}^{p} w_j^\top x_j$ with $x \in \mathbb{R}^p$ and $w_j \in \mathbb{R}$

  - by $f(x) = \sum_{j=1}^{p} w_j^\top \Phi_j(x)$ with $x \in \mathcal{X}$, $\Phi_j(x) \in \mathcal{F}_j$ an $w_j \in \mathcal{F}_j$

- Replace the $\ell_1$-norm $\sum_{j=1}^{p} |w_j|$ by "block" $\ell_1$-norm $\sum_{j=1}^{p} \|w_j\|_2$

- Remarks

  - Hilbert space extension of the group Lasso (Yuan and Lin, 2006)
  - Alternative sparsity-inducing norms (Ravikumar et al., 2008)

# Multiple kernel learning

- Learning combinations of kernels: $K(\eta) = \sum_{j=1}^{m} \eta_j K_j, \quad \eta \geqslant 0$

  - Summing kernels $\Leftrightarrow$ concatenating feature spaces
  - Assume $k_1(x, y) = \langle \Phi_1(x), \Phi_1(y) \rangle$, $k_2(x, y) = \langle \Phi_2(x), \Phi_2(y) \rangle$

$$k_1(x, y) + k_2(x, y) = \left\langle \begin{pmatrix} \Phi_1(x) \\ \Phi_2(x) \end{pmatrix}, \begin{pmatrix} \Phi_1(y) \\ \Phi_2(y) \end{pmatrix} \right\rangle$$

- Summing kernels $\Leftrightarrow$ generalized additive models

- Various priors on kernel weights $\eta$ do not change the function space

# Multiple kernel learning (MKL)
## (Lanckriet et al., 2004b; Bach et al., 2004a)

- Multiple feature maps / kernels on $x \in \mathcal{X}$:

  - $p$ "feature maps" $\Phi_j : \mathcal{X} \mapsto \mathcal{F}_j$, $j = 1, \ldots, p$.
  - Minimization with respect to $w_1 \in \mathcal{F}_1, \ldots, w_p \in \mathcal{F}_p$
  - Predictor: $f(x) = w_1^\top \Phi_1(x) + \cdots + w_p^\top \Phi_p(x)$

$$
\begin{array}{ccccc}
 & & \Phi_1(x)^\top & w_1 & \\
 & \nearrow & \vdots \quad \vdots & & \searrow \\
x & \longrightarrow & \Phi_j(x)^\top \quad w_j & \longrightarrow & w_1^\top \Phi_1(x) + \cdots + w_p^\top \Phi_p(x) \\
 & \searrow & \vdots \quad \vdots & & \nearrow \\
 & & \Phi_p(x)^\top & w_p &
\end{array}
$$

  - Generalized additive models (Hastie and Tibshirani, 1990)

# Regularization for multiple features

$$
\begin{array}{ccc}
 & \Phi_1(x)^\top & w_1 \\
\nearrow & \vdots \quad \vdots & \searrow \\
x \longrightarrow \ \Phi_j(x)^\top & w_j & \longrightarrow \ w_1^\top \Phi_1(x) + \cdots + w_p^\top \Phi_p(x) \\
\searrow & \vdots \quad \vdots & \nearrow \\
 & \Phi_p(x)^\top & w_p
\end{array}
$$

- Regularization by $\sum_{j=1}^{p} \|w_j\|_2^2$ is equivalent to using $K = \sum_{j=1}^{p} K_j$

  - Summing kernels is equivalent to concatenating feature spaces

# Regularization for multiple features

$$
\begin{array}{c}
\Phi_1(x)^\top \quad w_1 \\
\nearrow \qquad \vdots \qquad \vdots \qquad \searrow \\
x \longrightarrow \Phi_j(x)^\top \quad w_j \longrightarrow w_1^\top \Phi_1(x) + \cdots + w_p^\top \Phi_p(x) \\
\searrow \qquad \vdots \qquad \vdots \qquad \nearrow \\
\Phi_p(x)^\top \quad w_p
\end{array}
$$

- Regularization by $\sum_{j=1}^p \|w_j\|_2^2$ is equivalent to using $K = \sum_{j=1}^p K_j$

- Regularization by $\sum_{j=1}^p \|w_j\|_2$ imposes sparsity at the group level

- **Main results when regularizing by block $\ell_1$-norm**:

  1. Algorithms
  2. Analysis of sparsity inducing properties (Bach, 2008)
  3. Corresponds to a sparse combination of kernels
     $K = \sum_{j=1}^p \eta_j K_j$ (Bach et al., 2004a)

# Applications of multiple kernel learning

- **Selection of hyperparameters for kernel methods**

- **Fusion from heterogeneous data sources** (Lanckriet et al., 2004a)

- Two strategies for kernel combinations:

  - Uniform combination $\Leftrightarrow \ell_2$-norm
  - Sparse combination $\Leftrightarrow \ell_1$-norm
  - MKL always leads to more interpretable models
  - MKL does not always lead to better predictive performance
    - In particular, with few well-designed kernels
    - Be careful with normalization of kernels (Bach et al., 2004b)

- Rank-one kernel matrices: MKL $=$ Lasso

# Applications of multiple kernel learning

- **Selection of hyperparameters for kernel methods**

- **Fusion from heterogeneous data sources** (Lanckriet et al., 2004a)

- Two strategies for kernel combinations:

  - Uniform combination $\Leftrightarrow \ell_2$-norm
  - Sparse combination $\Leftrightarrow \ell_1$-norm
  - MKL always leads to more interpretable models
  - MKL does not always lead to better predictive performance
    - In particular, with few well-designed kernels
    - Be careful with normalization of kernels (Bach et al., 2004b)

- Rank-one kernel matrices: MKL $=$ Lasso

- **Sparse methods**: new possibilities and new features

# How does a good kernel matrix look like?

# How does a good kernel matrix look like?

# How does a good kernel matrix look like?

# How does a good kernel matrix look like?

# How does a good kernel matrix look like?

# How does a good kernel matrix look like?



- <span style="color:red">Good kernel matrices for classification may not be block-constants</span>

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Learning on matrices - Image denoising

- Simultaneously denoise all patches of a given image

- Example from Mairal, Bach, Ponce, Sapiro, and Zisserman (2009)

# Learning on matrices - Collaborative filtering

- Given $n_{\mathcal{X}}$ "movies" $\mathbf{x} \in \mathcal{X}$ and $n_{\mathcal{Y}}$ "customers" $\mathbf{y} \in \mathcal{Y}$,

- predict the "rating" $z(\mathbf{x}, \mathbf{y}) \in \mathcal{Z}$ of customer $\mathbf{y}$ for movie $\mathbf{x}$

- Training data: large $n_{\mathcal{X}} \times n_{\mathcal{Y}}$ incomplete matrix $\mathbf{Z}$ that describes the known ratings of some customers for some movies

- **Goal**: complete the matrix.

# Learning on matrices - Source separation

- Single microphone (Benaroya et al., 2006; Févotte et al., 2009)



Signal x

Log-power spectrogram

# Learning on matrices - Multi-task learning

- $k$ linear prediction tasks on same covariates $\mathbf{x} \in \mathbb{R}^p$

  - $k$ weight vectors $\mathbf{w}_j \in \mathbb{R}^p$
  - Joint matrix of predictors $\mathbf{W} = (\mathbf{w}_1, \ldots, \mathbf{w}_k) \in \mathbb{R}^{p \times k}$

- Classical application

  - Multi-category classification (one task per class) (Amit et al., 2007)

- **Share parameters between tasks**

- **Joint variable selection** (Obozinski et al., 2009)

  - Select variables which are predictive for all tasks

- **Joint feature selection** (Pontil et al., 2007)

  - Construct linear features common to all tasks

# Matrix factorization - Dimension reduction

- Given data matrix $\mathcal{X} = (\mathbf{x}_1, \ldots, \mathbf{x}_n) \in \mathbb{R}^{p \times n}$

  – **Principal component analysis**: $\boxed{\mathbf{x}_i \approx \mathbf{D}\alpha_i \Rightarrow \mathbf{X} = \mathbf{D}\mathbf{A}}$

  – **K-means**: $\boxed{\mathbf{x}_i \approx \mathbf{d}_k \Rightarrow \mathbf{X} = \mathbf{D}\mathbf{A}}$

# Two types of sparsity for matrices $M \in \mathbb{R}^{n \times p}$
## I - Directly on the elements of $M$

- Many zero elements: $M_{ij} = 0$



- Many zero rows (or columns): $(M_{i1}, \ldots, M_{ip}) = 0$

# Two types of sparsity for matrices $\mathbf{M} \in \mathbb{R}^{n \times p}$
## II - Through a factorization of $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$

- Matrix $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$, $\mathbf{U} \in \mathbb{R}^{n \times k}$ and $\mathbf{V} \in \mathbb{R}^{p \times k}$

- **Low rank**: $m$ small



- **Sparse decomposition**: $\mathbf{U}$ sparse

# Structured sparse matrix factorizations

- Matrix $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$, $\mathbf{U} \in \mathbb{R}^{n \times k}$ and $\mathbf{V} \in \mathbb{R}^{p \times k}$

- **Structure on $\mathbf{U}$ and/or $\mathbf{V}$**

  - Low-rank: $\mathbf{U}$ and $\mathbf{V}$ have few columns
  - Dictionary learning / sparse PCA: $\mathbf{U}$ has many zeros
  - Clustering ($k$-means): $\mathbf{U} \in \{0,1\}^{n \times m}$, $\mathbf{U}\mathbf{1} = \mathbf{1}$
  - Pointwise positivity: non negative matrix factorization (NMF)
  - Specific patterns of zeros (Jenatton et al., 2010)
  - Low-rank + sparse (Candès et al., 2009)
  - etc.

- **Many applications**

- **Many open questions** (Algorithms, identifiability, etc.)

# Multi-task learning

- Joint matrix of predictors $W = (w_1, \dots, w_k) \in \mathbb{R}^{p \times k}$

- **Joint variable selection** (Obozinski et al., 2009)

  - Penalize by the sum of the norms of rows of $W$ (group Lasso)
  - Select variables which are predictive for all tasks

# Multi-task learning

- Joint matrix of predictors $W = (w_1, \ldots, w_k) \in \mathbb{R}^{p \times k}$

- **Joint <span style="color:red">variable</span> selection** (Obozinski et al., 2009)

  - Penalize by the sum of the norms of rows of $W$ (group Lasso)
  - Select variables which are predictive for all tasks

- **Joint <span style="color:red">feature</span> selection** (Pontil et al., 2007)

  - Penalize by the trace-norm (see later)
  - Construct linear features common to all tasks

- Theory: allows number of observations which is sublinear in the number of tasks (Obozinski et al., 2008; Lounici et al., 2009)

- Practice: more interpretable models, slightly improved performance

# Low-rank matrix factorizations
## Trace norm

- Given a matrix $\mathbf{M} \in \mathbb{R}^{n \times p}$

  - Rank of $\mathbf{M}$ is the minimum size $m$ of **all** factorizations of $\mathbf{M}$ into $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$, $\mathbf{U} \in \mathbb{R}^{n \times m}$ and $\mathbf{V} \in \mathbb{R}^{p \times m}$
  - Singular value decomposition: $\mathbf{M} = \mathbf{U}\operatorname{Diag}(\mathbf{s})\mathbf{V}^\top$ where $\mathbf{U}$ and $\mathbf{V}$ have orthonormal columns and $\mathbf{s} \in \mathbb{R}^m_+$ are singular values

- Rank of $\mathbf{M}$ equal to the number of non-zero singular values

# Low-rank matrix factorizations
## Trace norm

- Given a matrix $\mathbf{M} \in \mathbb{R}^{n \times p}$

  - Rank of $\mathbf{M}$ is the minimum size $m$ of **all** factorizations of $\mathbf{M}$ into $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$, $\mathbf{U} \in \mathbb{R}^{n \times m}$ and $\mathbf{V} \in \mathbb{R}^{p \times m}$
  - Singular value decomposition: $\mathbf{M} = \mathbf{U}\operatorname{Diag}(\mathbf{s})\mathbf{V}^\top$ where $\mathbf{U}$ and $\mathbf{V}$ have orthonormal columns and $\mathbf{s} \in \mathbb{R}_+^m$ are singular values

- Rank of $\mathbf{M}$ equal to the number of non-zero singular values

- **Trace-norm (a.k.a. nuclear norm)** = sum of singular values

- Convex function, leads to a semi-definite program (Fazel et al., 2001)

- First used for collaborative filtering (Srebro et al., 2005)

- Multi-category classif. (Amit et al., 2007; Harchaoui et al., 2012)

# Sparse principal component analysis

- Given data $\mathcal{X} = (\mathbf{x}_1^\top, \ldots, \mathbf{x}_n^\top) \in \mathbb{R}^{p \times n}$, two views of PCA:

  - **Analysis view**: find the projection $\mathbf{d} \in \mathbb{R}^p$ of maximum variance (with deflation to obtain more components)
  - **Synthesis view**: find the basis $\mathbf{d}_1, \ldots, \mathbf{d}_k$ such that all $\mathbf{x}_i$ have low reconstruction error when decomposed on this basis

- For regular PCA, the two views are equivalent

# Sparse principal component analysis

- Given data $\mathcal{X} = (\mathbf{x}_1^\top, \ldots, \mathbf{x}_n^\top) \in \mathbb{R}^{p \times n}$, two views of PCA:

  - **Analysis view**: find the projection $\mathbf{d} \in \mathbb{R}^p$ of maximum variance (with deflation to obtain more components)
  - **Synthesis view**: find the basis $\mathbf{d}_1, \ldots, \mathbf{d}_k$ such that all $\mathbf{x}_i$ have low reconstruction error when decomposed on this basis

- For regular PCA, the two views are equivalent

- **Sparse extensions**

  - Interpretability
  - High-dimensional inference
  - Two views are differents
    - For analysis view, see d'Aspremont, Bach, and El Ghaoui (2008)

# Sparse principal component analysis
## Synthesis view

- Find $\mathbf{d}_1, \ldots, \mathbf{d}_k \in \mathbb{R}^p$ **sparse** so that

$$\sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_i \in \mathbb{R}^m} \left\| \mathbf{x}_i - \sum_{j=1}^{k} (\boldsymbol{\alpha}_i)_j \mathbf{d}_j \right\|_2^2 = \sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_i \in \mathbb{R}^m} \left\| \mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i \right\|_2^2 \text{ is small}$$

  - Look for $\mathbf{A} = (\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_n) \in \mathbb{R}^{k \times n}$ and $\mathbf{D} = (\mathbf{d}_1, \ldots, \mathbf{d}_k) \in \mathbb{R}^{p \times k}$ such that $\mathbf{D}$ is sparse and $\|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2$ is small

# Sparse principal component analysis
## Synthesis view

- Find $\mathbf{d}_1, \ldots, \mathbf{d}_k \in \mathbb{R}^p$ **sparse** so that

$$\sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_i \in \mathbb{R}^m} \left\| \mathbf{x}_i - \sum_{j=1}^{k} (\boldsymbol{\alpha}_i)_j \mathbf{d}_j \right\|_2^2 = \sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_i \in \mathbb{R}^m} \left\| \mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i \right\|_2^2 \text{ is small}$$

  - Look for $\mathbf{A} = (\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_n) \in \mathbb{R}^{k \times n}$ and $\mathbf{D} = (\mathbf{d}_1, \ldots, \mathbf{d}_k) \in \mathbb{R}^{p \times k}$ such that $\mathbf{D}$ is sparse and $\|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2$ is small

- Sparse formulation (Witten et al., 2009; Bach et al., 2008)

  - Penalize/constrain $\mathbf{d}_j$ by the $\ell_1$-norm for sparsity
  - Penalize/constrain $\boldsymbol{\alpha}_i$ by the $\ell_2$-norm to avoid trivial solutions

$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \sum_{j=1}^{k} \|\mathbf{d}_j\|_1 \text{ s.t. } \forall i, \|\boldsymbol{\alpha}_i\|_2 \leqslant 1$$

# Sparse PCA vs. dictionary learning

- **Sparse PCA**: $\mathbf{x}_i \approx \mathbf{D}\boldsymbol{\alpha}_i$, $\mathbf{D}$ sparse

# Sparse PCA vs. dictionary learning

- **Sparse PCA**: $\mathbf{x}_i \approx \mathbf{D}\boldsymbol{\alpha}_i$, $\mathbf{D}$ sparse



- **Dictionary learning**: $\mathbf{x}_i \approx \mathbf{D}\boldsymbol{\alpha}_i$, $\boldsymbol{\alpha}_i$ sparse

# Structured matrix factorizations (Bach et al., 2008)

$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \sum_{j=1}^{k} \|\mathbf{d}_j\|_\star \ \text{s.t.} \ \forall i, \|\boldsymbol{\alpha}_i\|_\bullet \leqslant 1$$

$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \sum_{i=1}^{n} \|\boldsymbol{\alpha}_i\|_\bullet \ \text{s.t.} \ \forall j, \|\mathbf{d}_j\|_\star \leqslant 1$$

- Optimization by alternating minimization (non-convex)

- $\boldsymbol{\alpha}_i$ decomposition coefficients (or "code"), $\mathbf{d}_j$ dictionary elements

- Two related/equivalent problems:

  - **Sparse PCA** = **sparse dictionary** ($\ell_1$-norm on $\mathbf{d}_j$)
  - **Dictionary learning** = **sparse decompositions** ($\ell_1$-norm on $\boldsymbol{\alpha}_i$) (Olshausen and Field, 1997; Elad and Aharon, 2006; Lee et al., 2007)

# Dictionary learning for image denoising



$$\underbrace{\mathbf{x}}_{\text{measurements}} = \underbrace{\mathbf{y}}_{\text{original image}} + \underbrace{\varepsilon}_{\text{noise}}$$

# Sparse methods for machine learning
## Why use sparse methods?

- **Sparsity as a proxy to interpretability**

  – Structured sparsity (Jenatton et al., 2009)

- **Sparsity for high-dimensional inference**

  – Influence on feature design

- **Sparse methods are not limited to least-squares regression**

- **Faster training/testing**

- **Better predictive performance?**

  – Problems are sparse if you look at them the right way

# Conclusion - Interesting questions/issues

- **Implicit vs. explicit features**

  – Can we algorithmically achieve $\log p = O(n)$ with explicit unstructured features?

- **Norm design**

  – What type of behavior may be obtained with sparsity-inducing norms?

- **Overfitting convexity**

  – Do we actually need convexity for matrix factorization problems?

# Course outline

1. **Losses for particular machine learning tasks**

   - Classification, regression, etc...

2. **Regularization by Hilbertian norms (kernel methods)**

   - Kernels and representer theorem
   - Convex duality, optimization and algorithms
   - Kernel methods
   - Kernel design

3. **Regularization by sparsity-inducing norms**

   - $\ell_1$-norm regularization
   - Theoretical results
   - Multiple kernel learning
   - Learning on matrices

# Conclusion - Interesting problems
# Machine learning for computer vision

- **Kernel design for computer vision**

  - Benefits of "kernelizing" existing representations
  - Combining kernels

- **Sparsity and computer vision**

  - Going beyond image denoising

- **Large numbers of classes**

  - Theoretical and algorithmic challenges

- **Structured output**

- **Semi-supervised learning**

# References

Y. Amit, M. Fink, N. Srebro, and S. Ullman. Uncovering shared structures in multiclass classification. In *Proceedings of the 24th international conference on Machine Learning (ICML)*, 2007.

N. Aronszajn. Theory of reproducing kernels. *Trans. Am. Math. Soc.*, 68:337–404, 1950.

F. Bach. Self-concordant analysis for logistic regression. Technical Report 0910.4627, ArXiv, 2009.

F. Bach, J. Mairal, and J. Ponce. Convex sparse matrix factorizations. Technical Report 0812.1869, ArXiv, 2008.

F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. Technical Report 00613125, HAL, 2011.

F. R. Bach. Consistency of the group Lasso and multiple kernel learning. *Journal of Machine Learning Research*, pages 1179–1225, 2008.

F. R. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3:1–48, 2002.

F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2004a.

F. R. Bach, R. Thibaux, and M. I. Jordan. Computing regularization paths for learning multiple kernels. In *Advances in Neural Information Processing Systems 17*, 2004b.

F. R. Bach, D. Heckerman, and E. Horvitz. Considering cost asymmetry in learning classifiers. *Journal of Machine Learning Research*, 7:1713–1741, 2006.

O. Banerjee, L. El Ghaoui, and A. d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *The Journal of Machine Learning Research*, 9: 485–516, 2008.

A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.

L. Benaroya, F. Bimbot, and R. Gribonval. Audio source separation with a single sensor. *IEEE Transactions on Speech and Audio Processing*, 14(1):191, 2006.

D. Bertsekas. *Nonlinear programming*. Athena Scientific, 1995.

P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.

Karsten M. Borgwardt, Cheng Soon Ong, Stefan Schönauer, S. V. N. Vishwanathan, Alex J. Smola, and Hans-Peter Kriegel. Protein function prediction via graph kernels. *Bioinformatics*, 21, 2005.

L. Bottou and C. J. Lin. Support vector machine solvers. In *Large scale kernel machines*, 2007.

Léon Bottou and Olivier Bousquet. Learning using large datasets. In *Mining Massive DataSets for Security*, NATO ASI Workshop Series. IOS Press, Amsterdam, 2008. URL `http://leon.bottou.org/papers/bottou-bousquet-2008b`. to appear.

E.J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Arxiv preprint arXiv:0912.3599*, 2009.

O. Chapelle. Training a support vector machine in the primal. *Neural Computation*, 19(5):1155–1178, 2007.

Scott Shaobing Chen, David L. Donoho, and Michael A. Saunders. Atomic decomposition by basis

pursuit. *SIAM Rev.*, 43(1):129–159, 2001. ISSN 0036-1445.

A. d'Aspremont, F. Bach, and L. El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9:1269–1294, 2008.

B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Ann. Stat.*, 32:407, 2004.

M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Trans. Image Proc.*, 15(12):3736–3745, 2006.

M. Fazel, H. Hindi, and S.P. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *Proceedings of the American Control Conference*, volume 6, pages 4734–4739, 2001.

C. Févotte, N. Bertin, and J.-L. Durrieu. Nonnegative matrix factorization with the itakura-saito divergence. with application to music analysis. *Neural Computation*, 21(3), 2009.

S. Fine and K. Scheinberg. Efficient SVM training using low-rank kernel representations. *Journal of Machine Learning Research*, 2:243–264, 2001.

P. A. Flach. The geometry of ROC space: understanding machine learning metrics through ROC isometrics. In *International Conference on Machine Learning (ICML)*, 2003.

W. Fu. Penalized regressions: the bridge vs. the Lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416, 1998).

Thomas Gärtner, Peter A. Flach, and Stefan Wrobel. On graph kernels: Hardness results and efficient alternatives. In *COLT*, 2003.

K. Grauman and T. Darrell. The pyramid match kernel: Efficient learning with sets of features. *J. Mach. Learn. Res.*, 8:725–760, 2007. ISSN 1533-7928.

Z. Harchaoui and F. R. Bach. Image classification with segmentation graph kernels. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007.

Z. Harchaoui, M. Douze, M. Paulin, M. Dudik, and J. Malick. Large-scale classification with trace-norm regularization. In *Proc. CVPR*, 2012.

T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer-Verlag, 2001.

T. Hastie, S. Rosset, R. Tibshirani, and J. Zhu. The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 5:1391–1415, 2005.

T. J. Hastie and R. J. Tibshirani. *Generalized Additive Models*. Chapman & Hall, 1990.

David Haussler. Convolution kernels on discrete structures. Technical report, UCSC, 1999.

M. Hein and O. Bousquet. Hilbertian metrics and positive-definite kernels on probability measures. In *AISTATS*, 2004.

R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. In *Advances in Large Margin Classifiers*. MIT Press, Cambridge, MA, 2000.

R. Jenatton, J.Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. Technical report, arXiv:0904.3523, 2009.

R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for sparse hierarchical dictionary learning. In *Submitted to ICML*, 2010.

T. Joachims. Making large-scale support vector machine learning practical. In *Advances in kernel methods — Support Vector learning*. MIT Press, 1998.

T. Joachims. Optimizing search engines using clickthrough data. In *ACM SIGKDD Conference on*

*Knowledge Discovery and Data Mining (KDD)*, pages 133–142, 2002.

Hisashi Kashima, Koji Tsuda, and Akihiro Inokuchi. Kernels for graphs. In *Kernel Methods in Computational Biology*. MIT Press, 2004.

G. S. Kimeldorf and G. Wahba. Some results on Tchebycheffian spline functions. *J. Math. Anal. Applicat.*, 33:82–95, 1971.

G. R. G. Lanckriet, T. De Bie, N. Cristianini, M. I. Jordan, and W. S. Noble. A statistical framework for genomic data fusion. *Bioinf.*, 20:2626–2635, 2004a.

G. R. G. Lanckriet, N. Cristianini, L. El Ghaoui, P. Bartlett, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, 5:27–72, 2004b.

S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. CVPR*, 2006.

H. Lee, A. Battle, R. Raina, and A. Ng. Efficient sparse coding algorithms. In *NIPS*, 2007.

G. Loosli, S. Canu, S. Vishwanathan, A. Smola, and M. Chattopadhyay. Boîte à outils SVM simple et rapide. *Revue dIntelligence Artificielle*, 19(4-5):741–767, 2005.

K. Lounici. Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators. *Electronic Journal of Statistics*, 2, 2008.

K. Lounici, A.B. Tsybakov, M. Pontil, and S.A. van de Geer. Taking advantage of sparsity in multi-task learning. In *Conference on Computational Learning Theory (COLT)*, 2009.

J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *International Conference on Computer Vision (ICCV)*, 2009.

H. M. Markowitz. The optimization of a quadratic function subject to linear constraints. *Naval*

*Research Logistics Quarterly*, 3:111–133, 1956.

N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *Annals of statistics*, 34(3):1436, 2006.

N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics*, 37(1):246–270, 2008.

F. Meyer. Hierarchies of partitions and morphological segmentation. In *Scale-Space and Morphology in Computer Vision*. Springer-Verlag, 2001.

Y. Nesterov. *Introductory lectures on convex optimization: A basic course*. Kluwer Academic Pub, 2003.

Y. Nesterov. Gradient methods for minimizing composite objective function. *Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, Tech. Rep*, 76, 2007.

G. Obozinski, M.J. Wainwright, and M.I. Jordan. High-dimensional union support recovery in multivariate regression. In *Advances in Neural Information Processing Systems (NIPS)*, 2008.

G. Obozinski, B. Taskar, and M.I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, pages 1–22, 2009.

B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37:3311–3325, 1997.

J. Platt. Fast training of support vector machines using sequential minimal optimization. In *Advances in Kernel Methods: Support Vector Learning*, 1998.

M. Pontil, A. Argyriou, and T. Evgeniou. Multi-task feature learning. In *Advances in Neural Information Processing Systems*, 2007.

Jan Ramon and Thomas Gärtner. Expressivity versus efficiency of graph kernels. In *First International Workshop on Mining Graphs, Trees and Sequences*, 2003.

P. Ravikumar, H. Liu, J. Lafferty, and L. Wasserman. SpAM: Sparse additive models. In *Advances in Neural Information Processing Systems (NIPS)*, 2008.

S. Rosset and J. Zhu. Piecewise linear regularized solution paths. *Ann. Statist.*, 35(3):1012–1030, 2007.

B. Schölkopf and A. J. Smola. *Learning with Kernels*. MIT Press, 2001.

B. Schölkopf, J. C. Platt, J. S. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, 2001.

S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal estimated sub-gradient solver for SVM. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2007.

J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Camb. U. P., 2004.

N. Srebro, J. D. M. Rennie, and T. S. Jaakkola. Maximum-margin matrix factorization. In *Advances in Neural Information Processing Systems 17*, 2005.

B. Taskar. Structured prediction: A large margin approach. In *NIPS Tutorial*, 2005. URL `media.nips.cc/Conferences/2007/Tutorials/Slides/taskar-NIPS-07-tutorial.ppt`.

R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of The Royal Statistical Society Series B*, 58(1):267–288, 1996.

S. A. Van De Geer. High-dimensional generalized linear models and the Lasso. *Annals of Statistics*, 36 (2):614, 2008.

R. Vert and J.-P. Vert. Consistency and convergence rates of one-class svms and related algorithms.

*Journal of Machine Learning Research*, 7:817–854, 2006.

S. V. N. Vishwanathan, A. J. Smola, and M. Murty. Simplesvm. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2003.

G. Wahba. *Spline Models for Observational Data*. SIAM, 1990.

M. J. Wainwright. Sharp thresholds for noisy and high-dimensional recovery of sparsity using $\ell_1$-constrained quadratic programming. Technical Report 709, Dpt. of Statistics, UC Berkeley, 2006.

C. Watkins. Dynamic alignment kernels. Technical report, RHUL, 1999.

D.M. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009.

Tong Tong Wu and Kenneth Lange. Coordinate descent algorithms for lasso penalized regression. *Ann. Appl. Stat.*, 2(1):224–244, 2008.

M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of The Royal Statistical Society Series B*, 68(1):49–67, 2006.

M. Yuan and Y. Lin. On the non-negative garrotte estimator. *Journal of The Royal Statistical Society Series B*, 69(2):143–161, 2007.

T. Zhang. Adaptive forward-backward greedy algorithm for sparse learning with linear models. *Advances in Neural Information Processing Systems*, 22, 2008.

T. Zhang. On the consistency of feature selection using greedy least squares regression. *The Journal of Machine Learning Research*, 10:555–568, 2009.

P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2563, 2006.

H. Zou. The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429, 2006.