

Machine Learning & Pattern Recognition

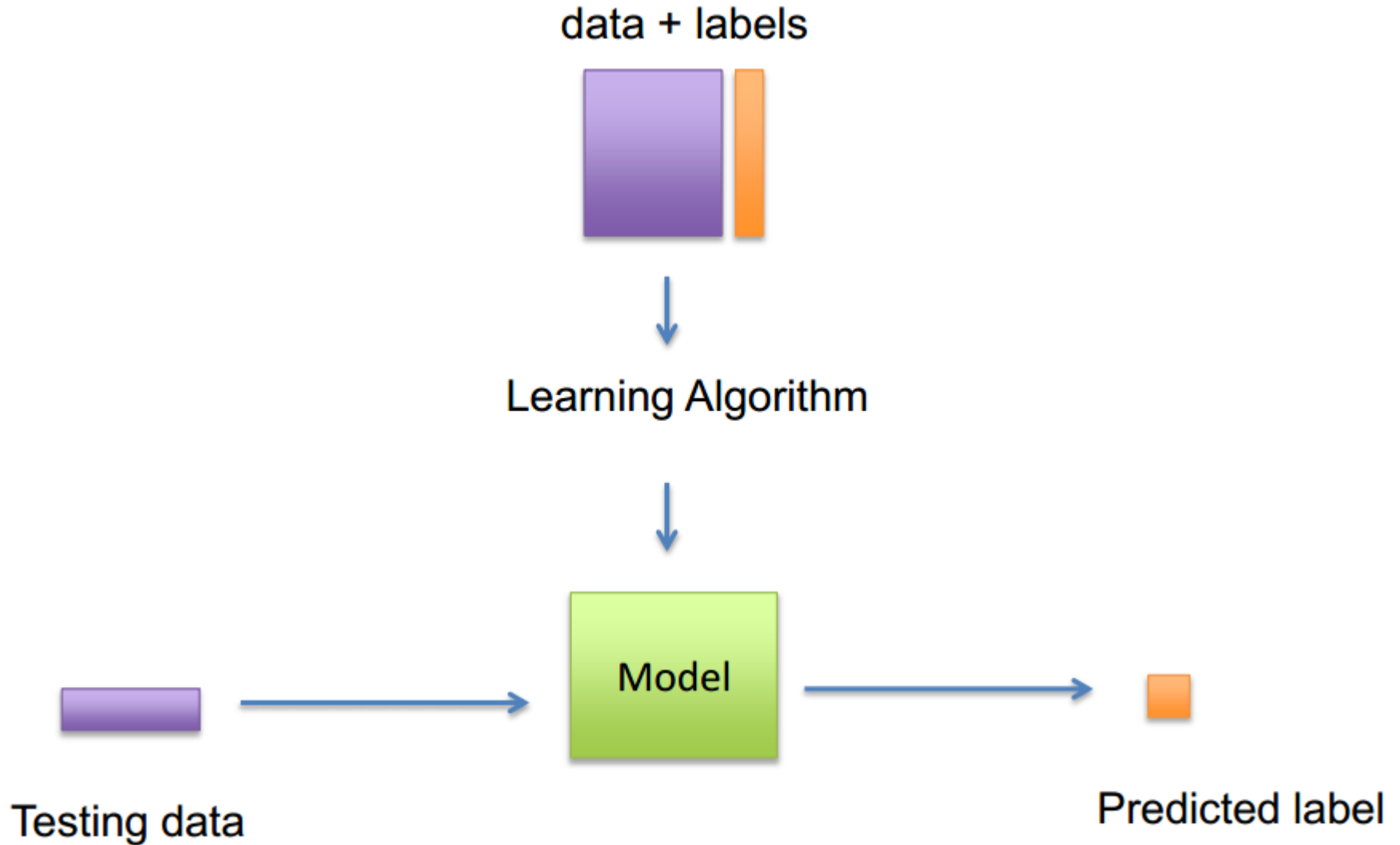
SONG Xuemeng

songxuемeng@sdu.edu.cn

<http://xuемeng.bitcron.com/>

Feature Selection

Supervised Learning



Text Classification

Is the story “interesting”?

It was a bright cold day in April, and the clocks were striking thirteen. Winston Smith, his chin nuzzled into his breast in an effort to escape the vile wind, slipped quickly through the glass doors of Victory Mansions, though not quickly enough to prevent a swirl of gritty dust from entering along with him.

Bag-of-words representation:

$x = \{0, 2, 0, 0, 1, \dots, 4, 0, 0, 0, 1\}$

One entry per word!

Easily 50,000 words! With big data...

- Time complexity
- Computational cost
- Overfitting
- Lack of interpretability

Feature selection

Some Things Matter, Some Do Not

- ***Relevant*** features
 - Those that we **need** to perform well
- ***Irrelevant*** features
 - Those that are simply **unnecessary**
- ***Redundant*** features
 - Those that **become** irrelevant in the presence of others

Feature Selection

Given: a set of features $X = \{X_1, X_2, \dots, X_D\}$ and a target variable Y .

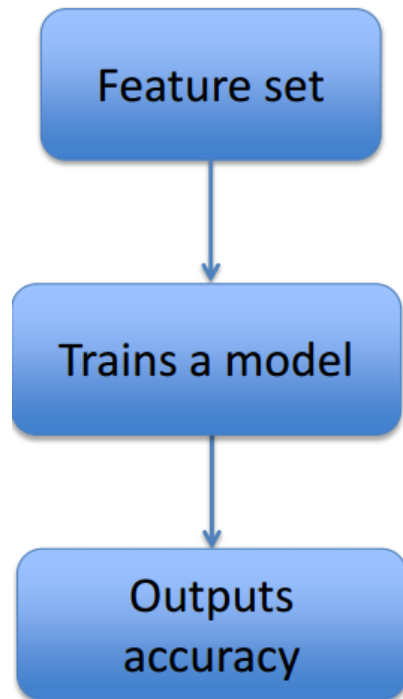
Find: minimum set S that achieves maximum classification performance of Y (for a given set of classifiers and classification performance metrics)

Feature Selection Techniques

- **Wrappers methods**
- **Filters methods**
- **Embedded methods**

Feature Selection (1): Wrapper Methods

Principle: We want to predict Y given the smallest possible subset of $\mathbf{X} = \{X_1, X_2, \dots, X_D\}$ while achieving maximal performance (e.g., accuracy).



- **Pros:**
 - Model-oriented
 - Usually gets good performance for the model you choose
- **Cons:**
 - Hugely computationally expensive

Feature Section (1): Wrapper Methods

With an **exhaustive** search

- With M features, 2^M possible feature subsets.

101110000001000100001000000000100101010

- 20 features... 1 million feature sets
- 25 features... 33.5 million sets
- 30 features... 1.1 billion sets

Feature Selection (1): Wrapper Methods

With an **exhaustive** search

- With M features, 2^M possible feature subsets.

101110000001000100001000000000100101010

- 20 features... 1 million feature sets
- 25 features... 33.5 million sets
- 30 features...1.1 billion sets

Need for a heuristic search strategy

1. Sequential forward selection

- ❖ Keep adding features one at a time until no further improvement can be achieved

2. Recursive backward elimination

- ❖ Start with the full set of predictors and keep removing features one at a time until no further improvement can be achieved

Wrappers: Sequential Forward Selection

Start with the empty set $S = \emptyset$

While stopping criteria not met

For each feature X_f not in S

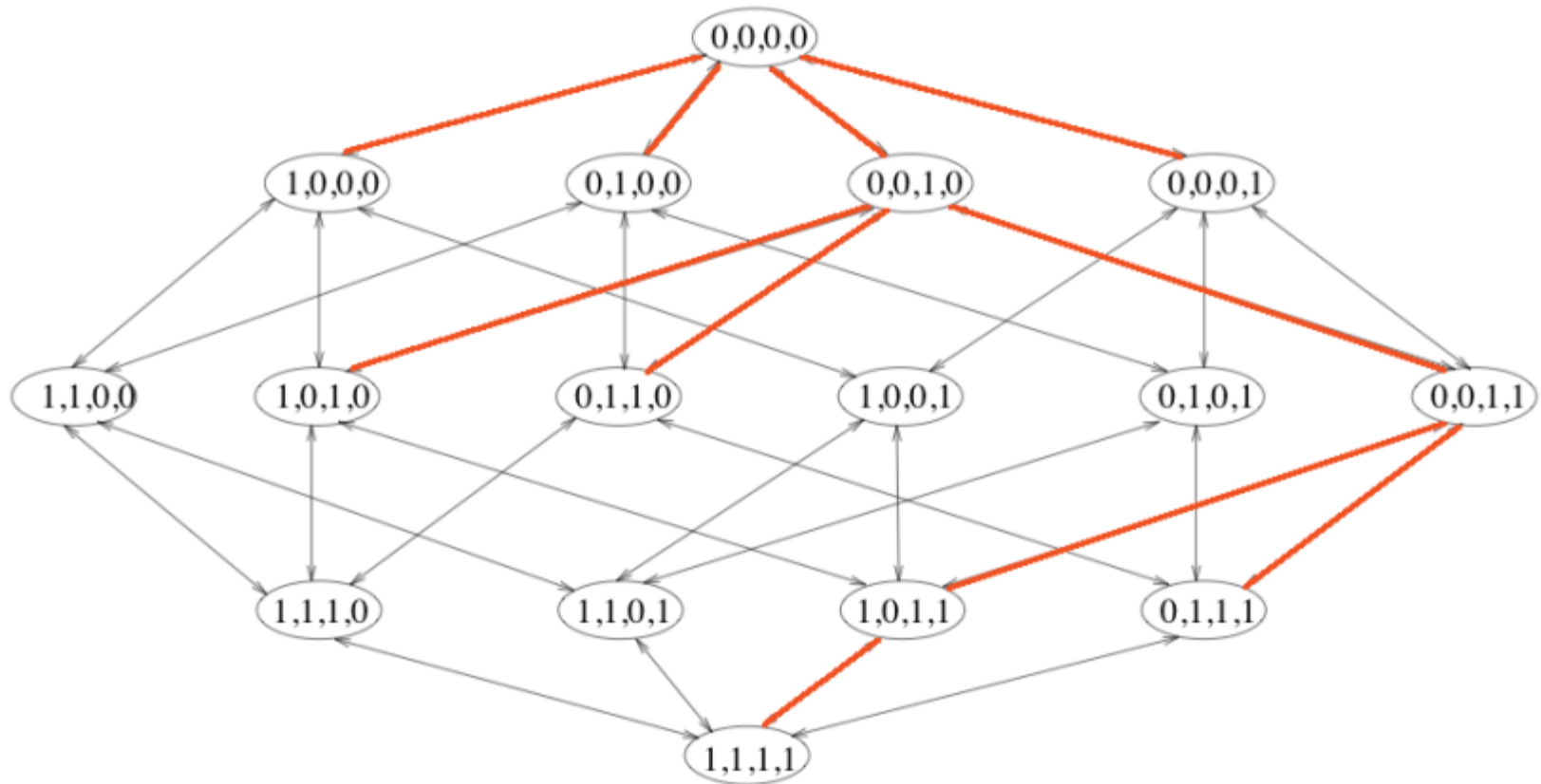
- Define $S' = S \cup \{X_f\}$
- Train model using the features in S'
- Compute the testing accuracy

End

$S = S'$ where S' is the feature set with the greatest accuracy

End

Search Complexity for Sequential Forward Selection



Evaluates $\frac{M(M+1)}{2}$ feature sets instead of 2^M !

Feature Selection (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05

1. Score each feature X_f individually based on the f-th column of the data matrix and label vector Y .

For each feature X_f

 Compute $J(X_f)$

End

2. Rank features according to $J(X_f)$.
3. Choose the top k features with the highest scores.

Feature Section (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

Examples of filtering criterion $J(X_f)$

- The mutual information
- Pearson r
- χ^2 -statistic
- Gini index

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05

Feature Selection (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

Examples of filtering criterion $J(X_f)$

- The mutual information

A measure of the mutual dependence between the two variables.

Score X_f based on the MI with Y .

$$J(X_f) = I(X_f, Y)$$

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05

Feature Section (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

Examples of filtering criterion $J(X_f)$

- Pearson r

A measure of the linear correlation between two variables X and Y .

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05

$$J(X_k) = \frac{\text{cov}(X_k, Y)}{\sqrt{\text{var}(X_k)}\sqrt{\text{var}(Y)}} \approx \frac{\sum_{i=1}^N (x_k^{(i)} - \bar{x}_k)(y^{(i)} - \bar{y})}{\sqrt{\sum_{i=1}^N (x_k^{(i)} - \bar{x}_k)^2} \sqrt{\sum_{i=1}^N (y^{(i)} - \bar{y})^2}}$$

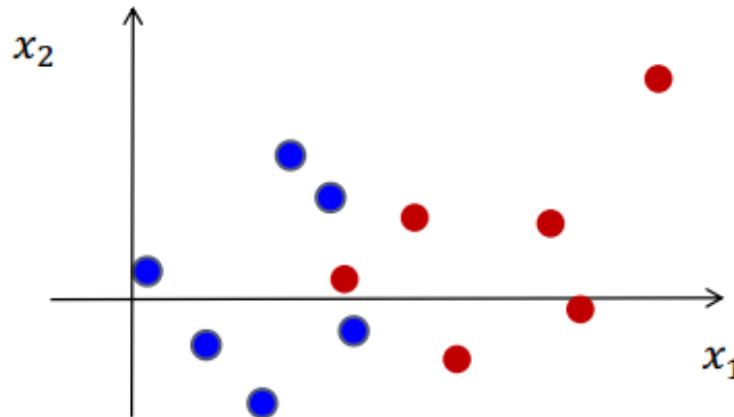
Feature Section (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

Examples of filtering criterion $J(X_f)$

- Pearson r

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05



$$J(x_1) \quad ? \quad J(x_2)$$

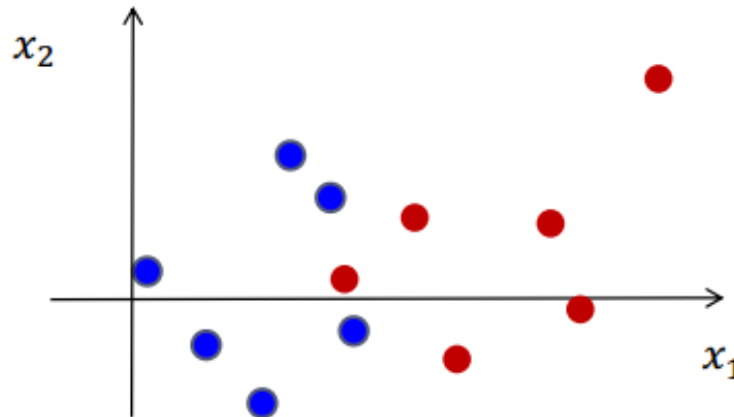
Feature Section (2): Filter Methods

Principle: replace evaluation of model with quick to compute statistics $J(X_f)$

Examples of filtering criterion $J(X_f)$

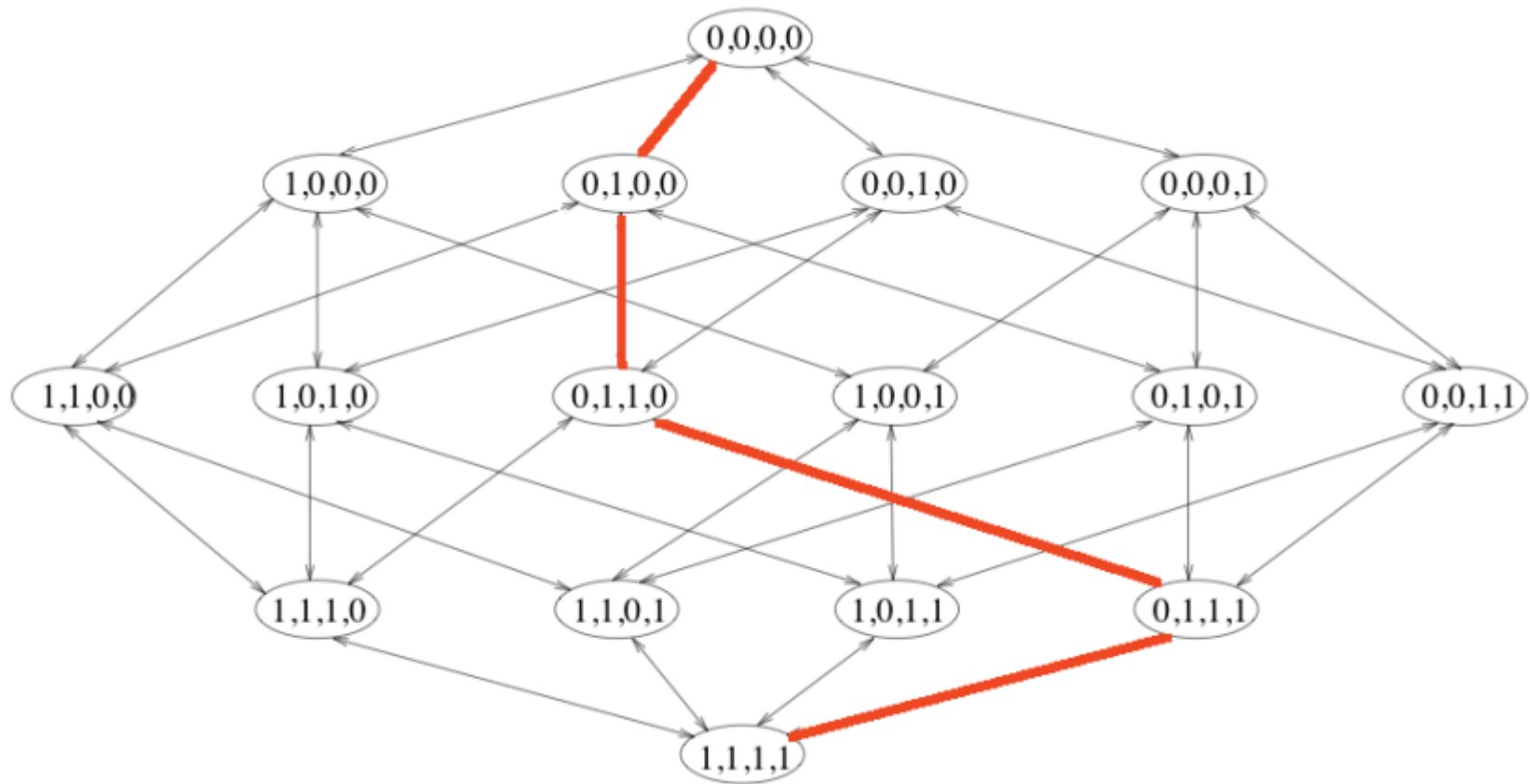
- Pearson r

k	$J(X_k)$
35	0.846
42	0.811
10	0.810
654	0.611
22	0.443
59	0.388
...	...
212	0.09
39	0.05



$$J(x_1) > J(x_2)$$

Search Complexity for Filter Methods



Pros:

- A lot less expensive!

Cons:

- Not model-oriented

Feature Section (3): Embedded methods

Principle: the classifier performs feature selection as part of the learning procedure.

Example: add the regularization term in the objective function

$$E_D(w) + \lambda E_W(w)$$

Underfitting & Overfitting

- The central challenge in machine learning is that we must perform well on **new, previously unseen** inputs—not just those on which our model was trained.
- The ability to perform well on previously unobserved inputs is called **generalization**.
- We typically estimate the **generalization error** (also called **test error**) of a model by measuring its performance on a *test set* which is collected separately from the training set.

Underfitting & Overfitting

- **Typically,**
 1. We sample the training set
 2. Then use it to choose the parameters to reduce training set error
 3. Then evaluate the model with the test set.
 - **Under this process,**
 - the expected test error \geq the training error.
-
- The factors determining how well a model will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.

Underfitting & Overfitting

- The factors determining how well a model will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.
- **Underfitting** occurs when the model is not able to obtain a sufficiently low error value on the **training** set.
- **Overfitting** occurs when the **gap** between the **training** error and **test** error is too large.

Underfitting & Overfitting

- We can control whether a model is more likely to overfit or underfit by altering its **capacity**.
- Informally, a model's capacity is its ability to fit a wide variety of functions.
- One way to control the capacity of a learning algorithm is by choosing its **hypothesis space**, the set of functions that the learning algorithm is allowed to select as being the solution.

Linear regression

$$y = b + wx$$

Introduce x^2 (quadratic model)

$$y = b + w_1x + w_2x^2$$

Continue to add more powers of x

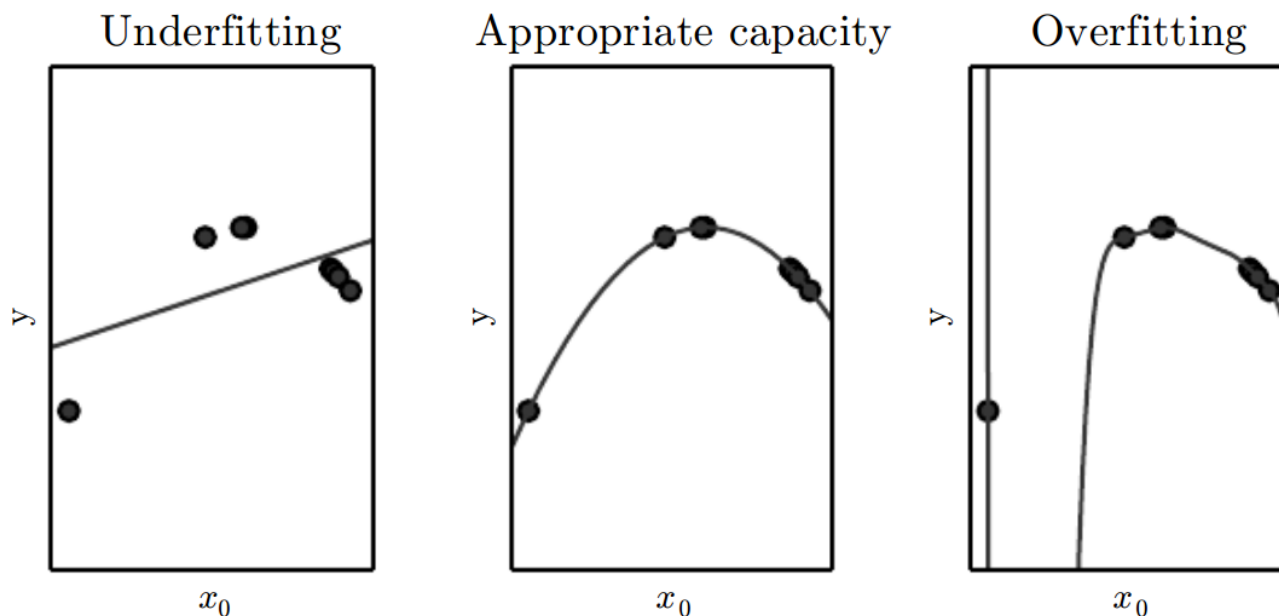
$$y = b + \sum_{i=1}^9 w_i x^i$$

Underfitting & Overfitting

Machine learning algorithms will generally perform best when their capacity is appropriate in regard to the **true complexity** of the task they need to perform and the **amount of training data** they are provided with.

Underfitting & Overfitting

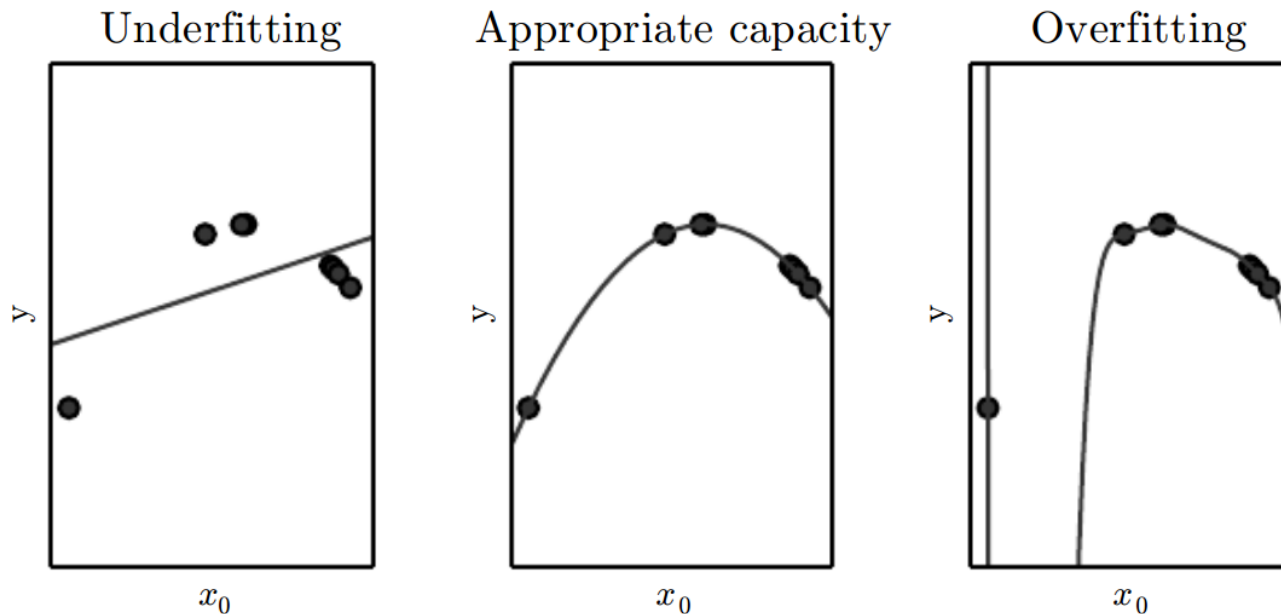
Machine learning algorithms will generally perform best when their capacity is appropriate in regard to the **true complexity** of the task they need to perform and the **amount of training data** they are provided with.



We compare a linear, quadratic and degree-9 predictor attempting to fit a problem where the **true** underlying function is quadratic.

Underfitting & Overfitting

- Models with low capacity may struggle to fit the training set.
- Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.

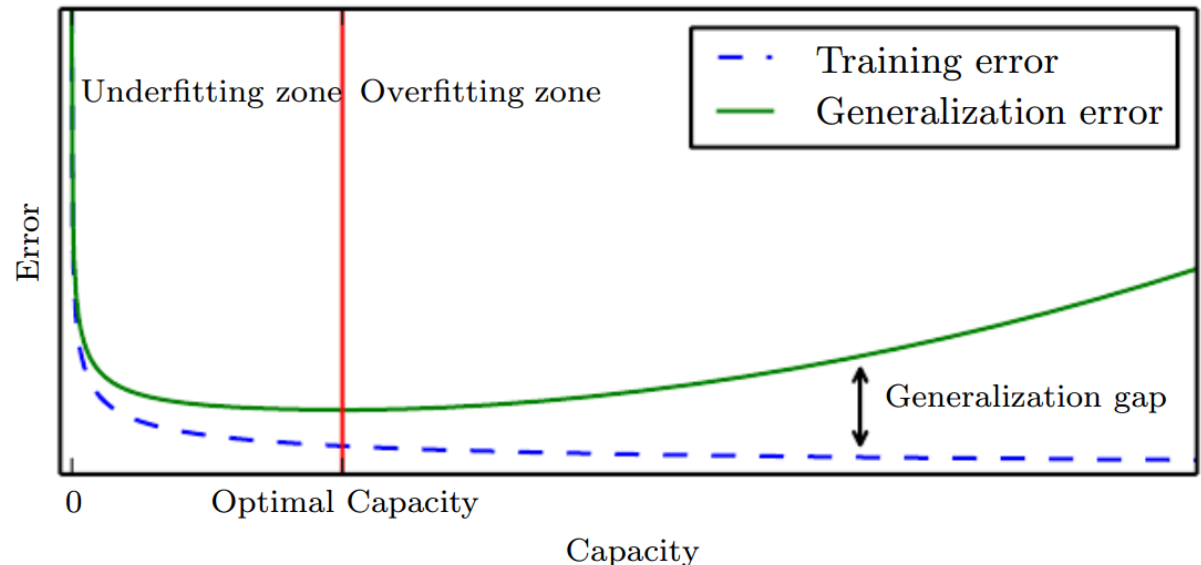


We compare a linear, quadratic and degree-9 predictor attempting to fit a problem where the **true** underlying function is quadratic.

Underfitting & Overfitting

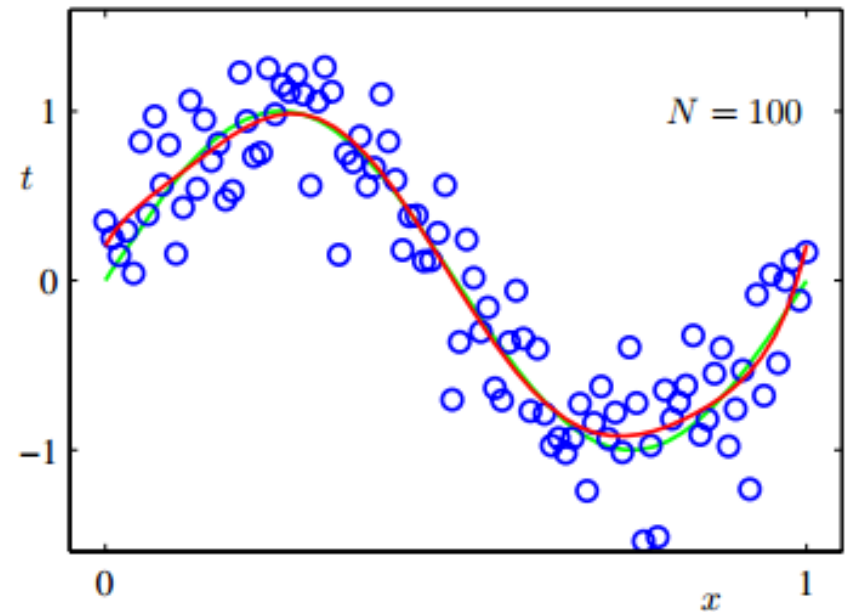
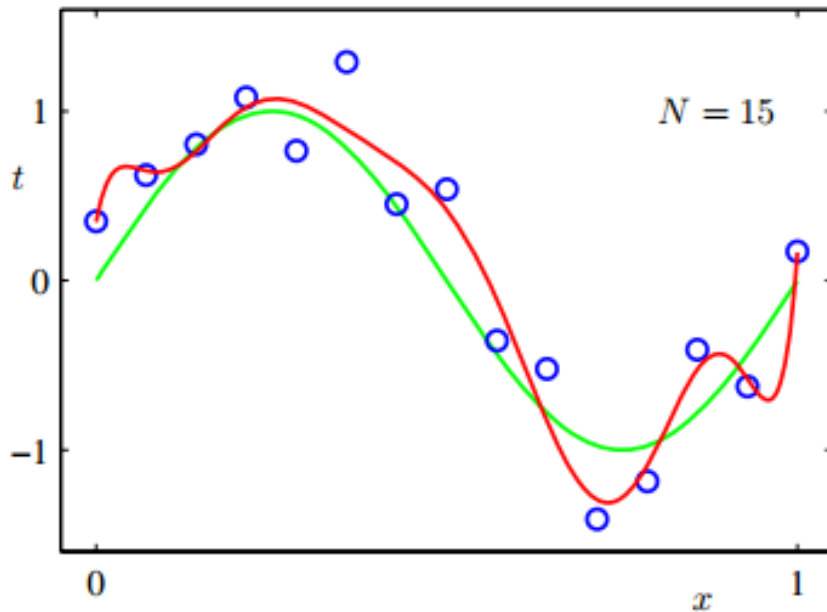
- Simpler functions are more likely to generalize.
 - Also need to choose a sufficiently complex hypothesis to achieve low training error.
- Training error decreases until it asymptotes to the minimum possible error value as model capacity increases.
 - Generalization error has a U-shaped curve as a function of model capacity.

Typical relationship between capacity and error.



Underfitting & Overfitting

- Increasing the size of dataset reduces the over-fitting problem.
- The larger the data set, the more complex the model that we can afford to fit to the data.



Plots of the solutions obtained by minimizing the sum-of-squares error function using the $M = 9$ polynomial for $N = 15$ data points (left plot) and $N = 100$ data points (right plot).

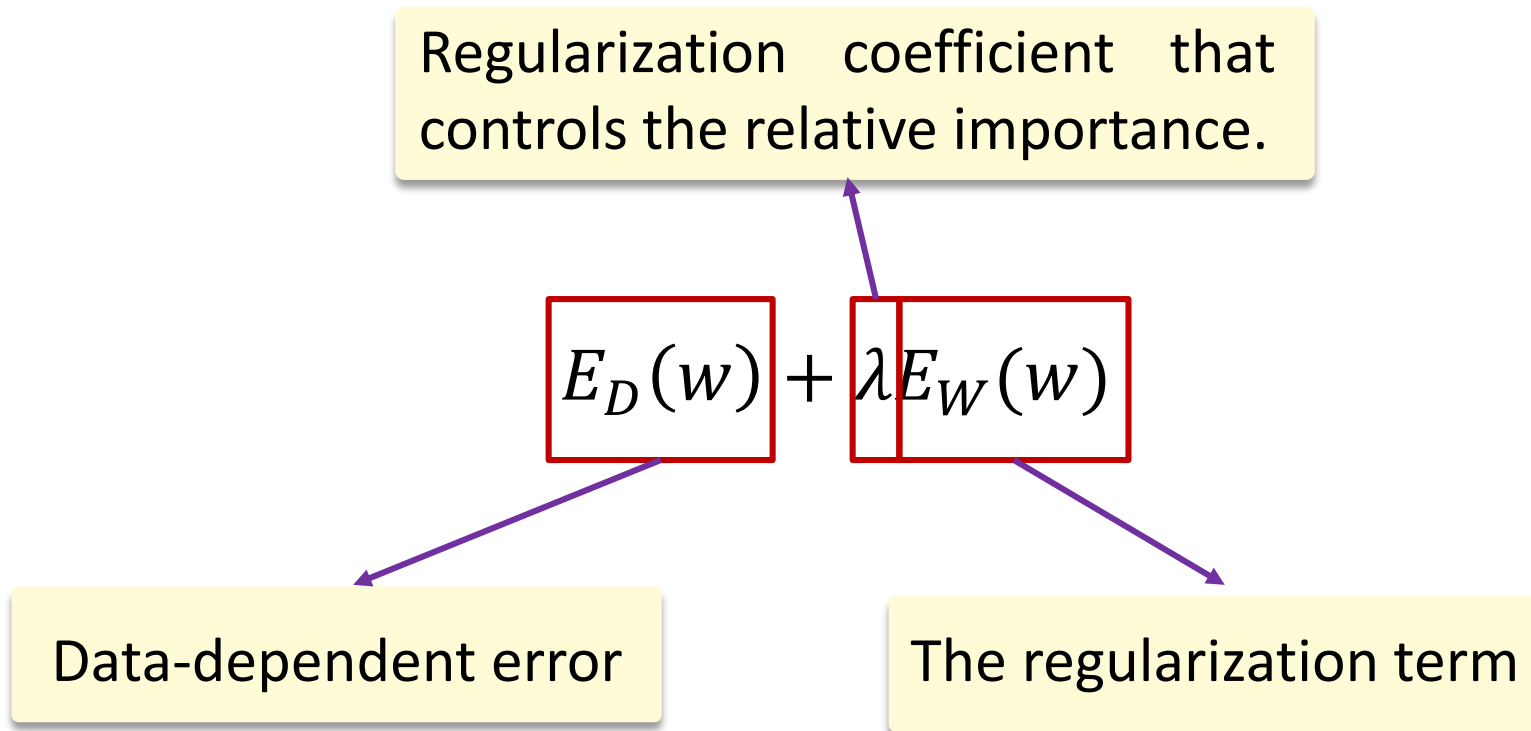
Avoid Overfitting

Adding a regularization term

$$E_D(w) + \lambda E_W(w)$$

Avoid Overfitting

Adding a regularization term



Avoid Overfitting

Adding a regularization term

Regularization coefficient that controls the relative importance.

$$E_D(w) + \lambda E_W(w)$$

Data-dependent error

The regularization term

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

We simply define
 $\phi(\mathbf{x}) = \mathbf{x}$

Avoid Overfitting

Adding a regularization term

$$E_D(w) + \lambda E_W(w)$$

One simple form of regularizer

$$E_W(w) = \frac{1}{2} w^T w \qquad E_T(w) = E_D(w) + \frac{\lambda}{2} w^T w$$

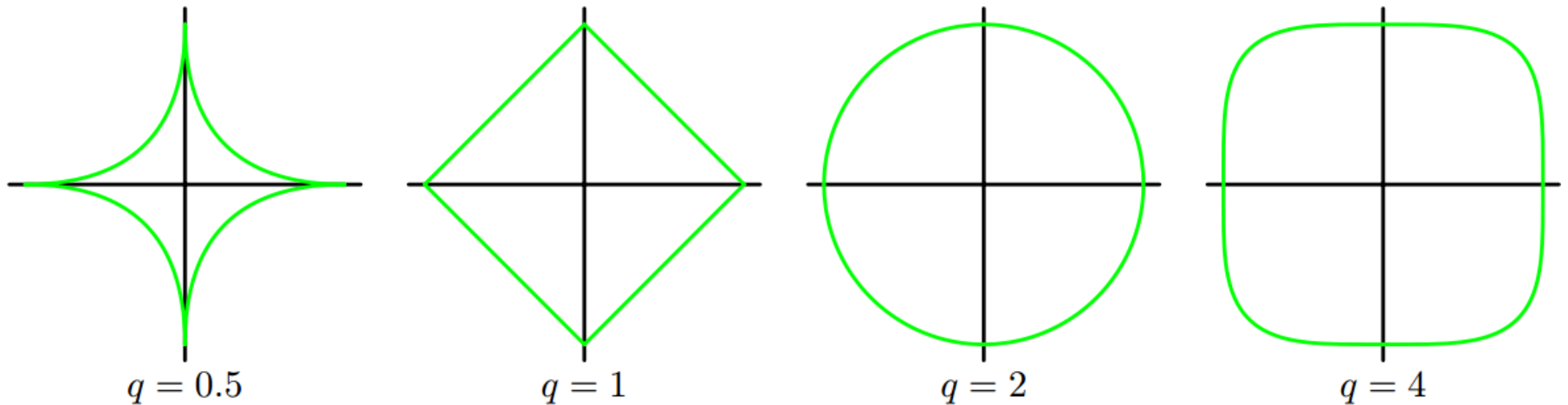
weight decay: encourages weight values to decay towards zero.

parameter shrinkage: shrinks parameter values towards zero.

❖ Advantage: Remains a quadratic function of w , so its exact minimizer can be found in closed form.

More General Regularizer

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

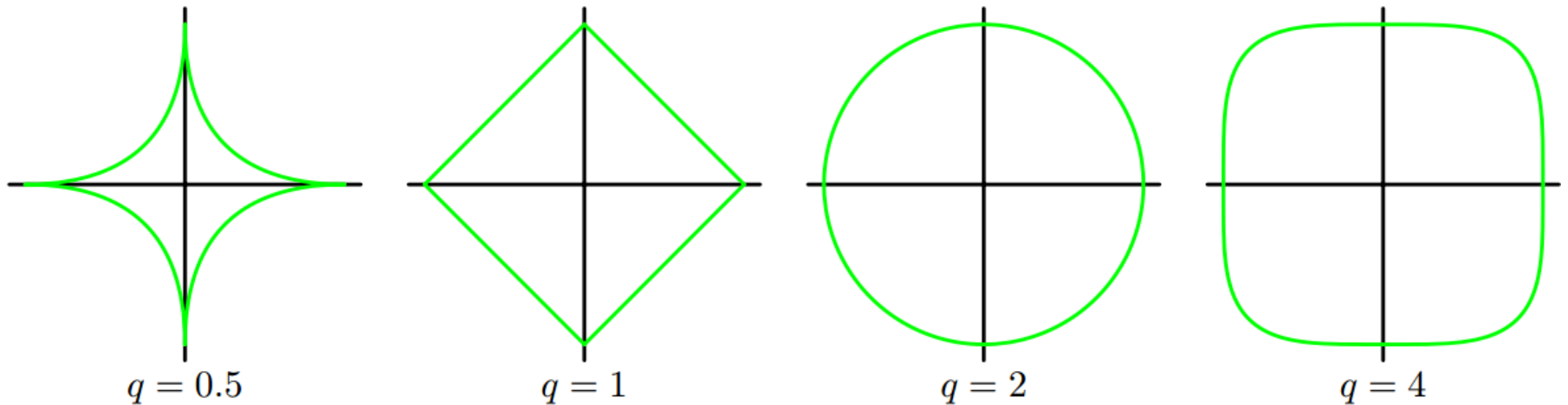


Contours of the regularization term for various values of the q .

- $q = 2$ corresponds to the *weight decay*.
- $q = 1$ corresponds to the *lasso*.
 - If λ is sufficiently large, some of the coefficients w_j are driven to zero, leading to a *sparse* model.

More General Regularizer

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



Contours of the regularization term for various values of the q .

- $q = 2$ corresponds to the *weight decay*.
- $q = 1$ corresponds to the *lasso*.
 - If λ is sufficiently large, some of the coefficients w_j are driven to zero, leading to a *sparse* model.

WHY?

More General Regularizer

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

Note that minimizing the above function is equivalent to minimizing the **unregularized** sum-of-squares error $E_D(w)$ subject to the constraint

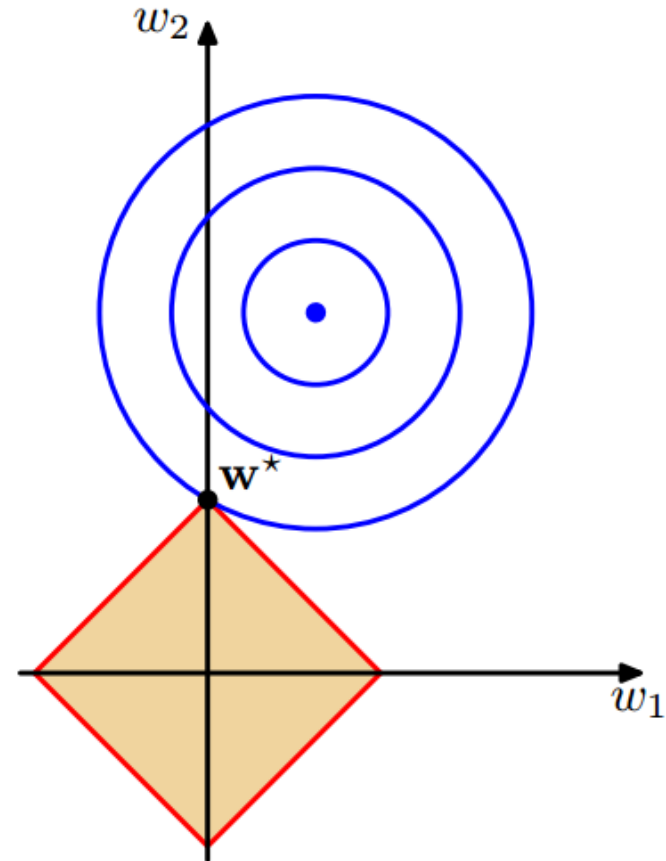
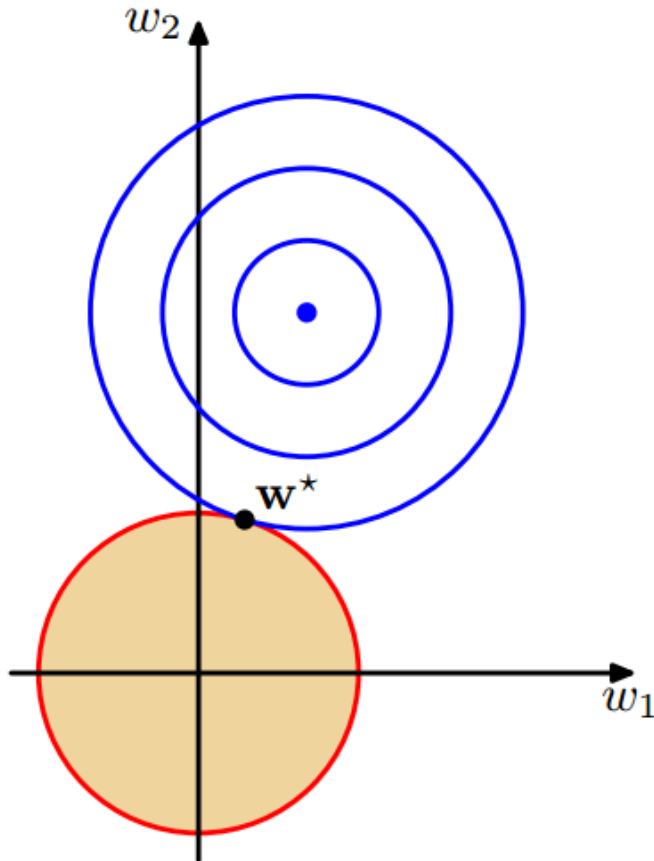
$$\sum_{j=1}^M |w_j|^q \leq \eta$$

for an appropriate value of the parameter η (*Lagrange multipliers*).

L1 VS L2 Regularization

w^* is the optimum value for w .

The lasso gives a sparse solution ($w_1^* = 0$).

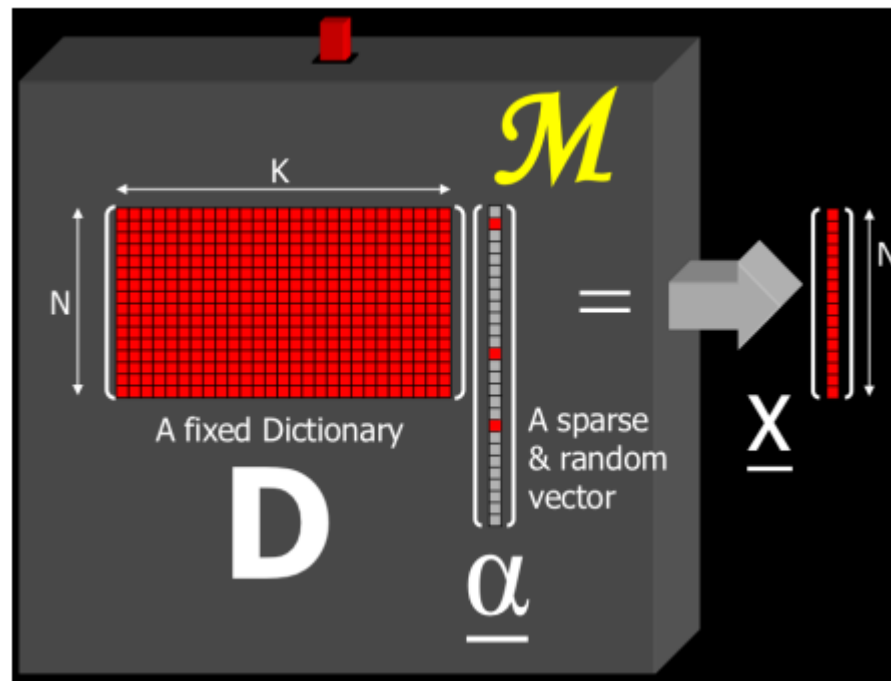


The contours of the unregularized error function (blue) along with the constraint region for the weight decay $q = 2$ (left) and the lasso $q = 1$ (right).

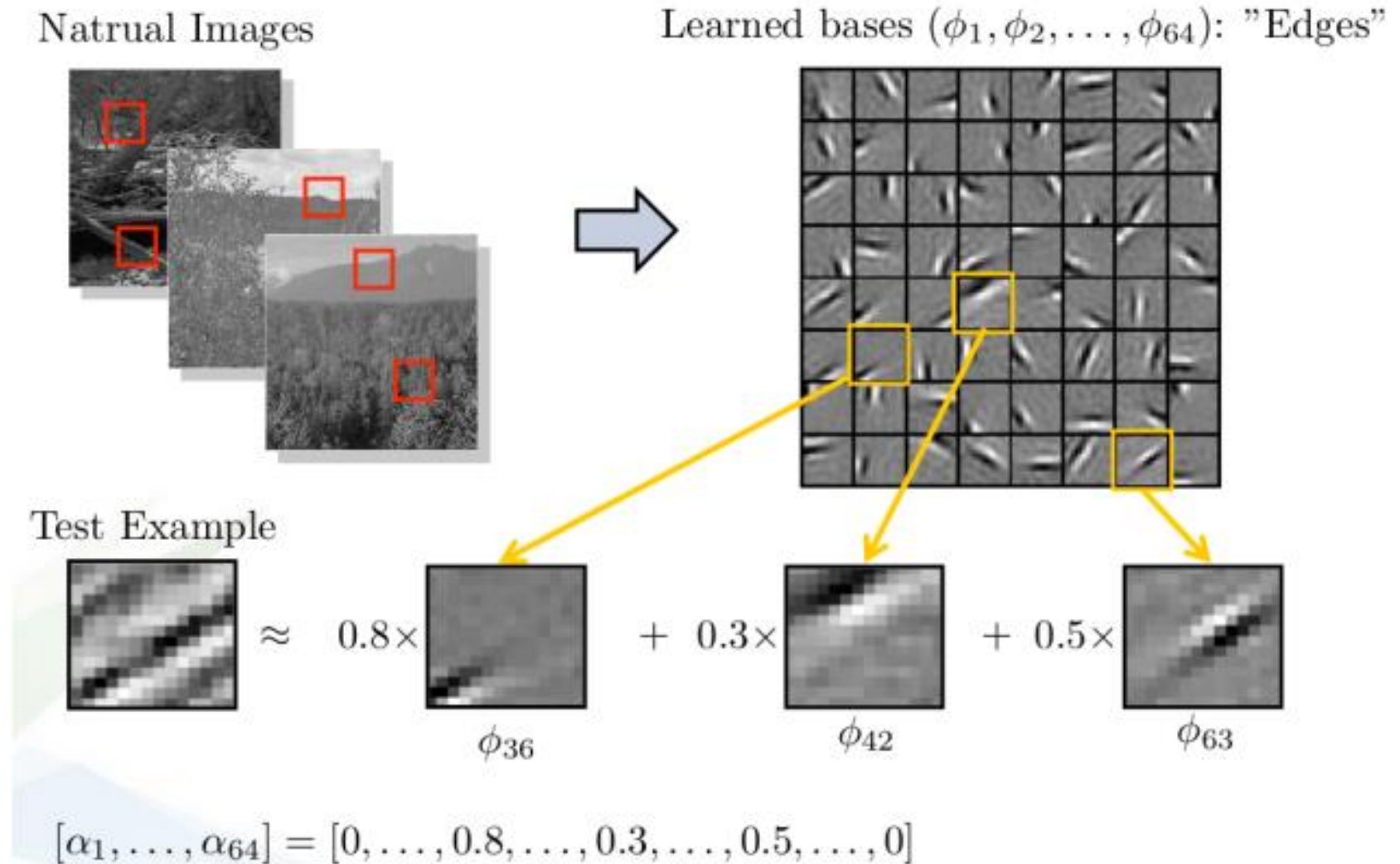
Sparse Coding

Every column of \mathbf{D} is a prototype (basis).

Similar to, but more general than, PCA.



Example: Sparse Coding---Images



Example: Sparse Coding---NLP

a	1	\approx	$0.7 \times$	ϕ_{36}	$0.4 \times$	ϕ_{42}	$0.1 \times$	ϕ_{63}			
aardvark	0								.1	0	.1
aardwolf	0								.2	.1	0
able	0								0	0	0.8
account	1								0	0	0
acid	0								.1	.1	0
across	0								.1	0	.1
...	0								0	0	.1
baby	1								0	.2	.2
back	0								1.2	0	.1
...	0								0	0	0
cradle	0								0	.2	0
...	1								0	.1	.3
...	0								1	0	0
...	0								0	.1	0
...	0	.1	.1	0							
zylophone	1	0	0	1.2							

Document $\approx 0.7 \times \text{Topic36} + 0.4 \times \text{Topic42} + 0.1 \times \text{Topic63}$

Sparse Coding

We assume our data \mathbf{x} satisfies

$$\mathbf{x} = \sum_{i=0}^k \alpha_i \boldsymbol{\phi}_i = \boldsymbol{\alpha} \mathbf{D}$$

Learning:

Given training data $\mathbf{x}^j, j \in \{1, \dots, m\}$

Learn dictionary \mathbf{D} and **sparse** code $\boldsymbol{\alpha}$ (i.e., α_i are mostly zero)

Encoding:

Given test data \mathbf{x} , learned dictionary \mathbf{D}

Compute the sparse code $\boldsymbol{\alpha}$

Sparse Coding

Input: Images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m$ (each in $\mathbb{R}^{n \times n}$)

Learn: Dictionary of bases $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{\alpha, \phi} \sum_{i=1}^m \left(\left\| \mathbf{x}^i - \sum_{j=1}^k \alpha_j^i \boldsymbol{\phi}_j \right\|^2 + \lambda \sum_{j=1}^k |\alpha_j^i| \right)$$

L1 sparsity term

Causes most α_j^i s to be 0.

Sparse Coding

Input: Images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m$ (each in $\mathbb{R}^{n \times n}$)

Learn: Dictionary of bases $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{\alpha, \phi} \sum_{i=1}^m \left(\left\| \mathbf{x}^i - \sum_{j=1}^k \alpha_j^i \boldsymbol{\phi}_j \right\|^2 + \lambda \sum_{j=1}^k |\alpha_j^i| \right)$$

L1 sparsity term

Causes most α_j^i s to be 0.

How to optimize α and ϕ ?

Sparse Coding

Input: Images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m$ (each in $\mathbb{R}^{n \times n}$)

Learn: Dictionary of bases $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{\alpha, \phi} \sum_{i=1}^m \left(\left\| \mathbf{x}^i - \sum_{j=1}^k \alpha_j^i \boldsymbol{\phi}_j \right\|^2 + \lambda \sum_{j=1}^k |\alpha_j^i| \right)$$

L1 sparsity term

Causes most α_j^i 's to be 0.

Alternating minimization:

Alternately minimize with respect to $\boldsymbol{\phi}_j$'s (easy) and α_j 's (harder).

Sparse Coding

Training time

Input: Images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m$ (each in $\mathbb{R}^{n \times n}$)

Learn: Dictionary of bases $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{\alpha, \phi} \sum_{i=1}^m \left(\|\mathbf{x}^i - \sum_{j=1}^k \alpha_j^i \boldsymbol{\phi}_j\|^2 + \lambda \sum_{j=1}^k |\alpha_j^i| \right)$$

Input: New image \mathbf{x} (in $\mathbb{R}^{n \times n}$) and previously learned $\boldsymbol{\phi}_i$'s.

Output: Representation $[\alpha_1, \alpha_2, \dots, \alpha_n]$ of image \mathbf{x} .

Test time

$$\min_{\alpha} \sum_{i=1}^m \left(\|\mathbf{x} - \sum_{j=1}^k \alpha_j \boldsymbol{\phi}_j\|^2 + \lambda \sum_{j=1}^k |\alpha_j| \right)$$

Sparse Coding

Training time

Input: Images $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m$ (each in $\mathbb{R}^{n \times n}$)

Learn: Dictionary of bases $\phi_1, \phi_2, \dots, \phi_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{\alpha, \phi} \sum_{i=1}^m \left(\|\mathbf{x}^i - \sum_{j=1}^k \alpha_j^i \phi_j\|^2 + \lambda \sum_{j=1}^k |\alpha_j^i| \right)$$

Input: New image \mathbf{x} (in $\mathbb{R}^{n \times n}$) and previously learned ϕ_i 's.

Output: Representation $[\alpha_1, \alpha_2, \dots, \alpha_n]$ of image \mathbf{x} .

Test time


$$\mathbf{x} \approx 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{63}$$

Represent as: $[0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, \dots]$

Optimization with L1 Norm

A simple form of the problem with L1 norm is defined as follows,

$$\begin{aligned}\alpha^* &= \min_{\alpha} h(\alpha) = \lambda \|\alpha\|_1 + \frac{1}{2} \|\alpha - s\|^2 \\ &= \sum_{j=1}^d \lambda |\alpha_j| + \sum_{j=1}^d \frac{1}{2} (\alpha_j - s_j)^2\end{aligned}$$

where α^* is the optimal solution. We can conclude that,

(1) If $\alpha_j > 0$, then $h(\alpha_j) = \lambda \alpha_j + \frac{1}{2} (\alpha_j - s_j)^2$ and its derivative is $h'(\alpha_j) = \lambda + \alpha_j^* - s_j = 0 \Rightarrow \alpha_j^* = s_j - \lambda$, which indicates $s_j > \lambda$;

Optimization with L1 Norm

A simple form of the problem with L1 norm is defined as follows,

$$\begin{aligned}\alpha^* &= \min_{\alpha} h(\alpha) = \lambda \|\alpha\|_1 + \frac{1}{2} \|\alpha - s\|^2 \\ &= \sum_{j=1}^d \lambda |\alpha_j| + \sum_{j=1}^d \frac{1}{2} (\alpha_j - s_j)^2\end{aligned}$$

where α^* is the optimal solution. We can conclude that,

- (1) If $\alpha_j > 0$, then $h(\alpha_j) = \lambda \alpha_j + \frac{1}{2} (\alpha_j - s_j)^2$ and its derivative is $h'(\alpha_j) = \lambda + \alpha_j^* - s_j = 0 \Rightarrow \alpha_j^* = s_j - \lambda$, which indicates $s_j > \lambda$;
- (2) If $\alpha_j < 0$, then $h(\alpha_j) = -\lambda \alpha_j + \frac{1}{2} (\alpha_j - s_j)^2$ and its derivative is $h'(\alpha_j) = -\lambda + \alpha_j^* - s_j = 0 \Rightarrow \alpha_j^* = s_j + \lambda$, which indicates $s_j < -\lambda$;

Optimization with L1 Norm

A simple form of the problem with L1 norm is defined as follows,

$$\begin{aligned}\alpha^* &= \min_{\alpha} h(\alpha) = \lambda \|\alpha\|_1 + \frac{1}{2} \|\alpha - s\|^2 \\ &= \sum_{j=1}^d \lambda |\alpha_j| + \sum_{j=1}^d \frac{1}{2} (\alpha_j - s_j)^2\end{aligned}$$

where α^* is the optimal solution. We can conclude that,

- (1) If $\alpha_j > 0$, then $h(\alpha_j) = \lambda \alpha_j + \frac{1}{2} (\alpha_j - s_j)^2$ and its derivative is $h'(\alpha_j) = \lambda + \alpha_j^* - s_j = 0 \Rightarrow \alpha_j^* = s_j - \lambda$, which indicates $s_j > \lambda$;
- (2) If $\alpha_j < 0$, then $h(\alpha_j) = -\lambda \alpha_j + \frac{1}{2} (\alpha_j - s_j)^2$ and its derivative is $h'(\alpha_j) = -\lambda + \alpha_j^* - s_j = 0 \Rightarrow \alpha_j^* = s_j + \lambda$, which indicates $s_j < -\lambda$;
- (3) If $-\lambda \leq s_j \leq \lambda$, and then $\alpha_j^* = 0$.

Shrink Operator

α^* can be summarized as follows,

$$\alpha_j^* = \begin{cases} s_j - \lambda, & \text{if } s_j > \lambda \\ s_j + \lambda, & \text{if } s_j < -\lambda \\ 0 & \text{otherwise} \end{cases}$$

The equivalent expression is $\alpha^* = \textit{shrink}(\mathbf{s}, \lambda)$, where the j -th component of $\textit{shrink}(\mathbf{s}, \lambda)$,

$$\textit{shrink}(s, \lambda)_j = \textit{sign}(s_j) \max\{|s_j| - \lambda, 0\}$$

ISTA (Iterative Shrinkage-Thresholding Algorithms)

The objective function of ISTA has the form of

$$\arg \min F(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 = f(\boldsymbol{\alpha}) + g(\boldsymbol{\alpha})$$

which is usually difficult to solve. We convert it to the easy one.

First, the second order Taylor expansion is used to approximate

$f(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\alpha} - \mathbf{y}\|_2^2$ at point $\boldsymbol{\alpha}^{k-1}$:

$$f(\boldsymbol{\alpha}) = f(\boldsymbol{\alpha}^{k-1}) + (\boldsymbol{\alpha} - \boldsymbol{\alpha}^{k-1})^T \nabla f(\boldsymbol{\alpha}^{k-1}) + \frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^{k-1})^T \underline{\underline{H_f(\boldsymbol{\alpha}^{k-1})}} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^{k-1})$$

Hessian matrix of $f(\boldsymbol{\alpha})$ at $\boldsymbol{\alpha}^{k-1}$

ISTA

A real-valued function $f : \mathbb{R} \rightarrow \mathbb{R}$ is called **Lipschitz continuous** if there exists a **positive real constant** L such that, for all real x_1 and x_2 ,

$$|f(x_1) - f(x_2)| \leq L|x_1 - x_2| \quad L: \text{Lipschitz constant.}$$

Simple algebra shows that given any α^{k-1} and α , we always have,

$$f(\alpha) \leq f(\alpha^{k-1}) + (\alpha - \alpha^{k-1})^T \nabla f(\alpha^{k-1}) + \frac{L}{2} \|\alpha - \alpha^{k-1}\|^2$$

We thus approximate $\arg \min_{\alpha} F(\alpha) = \arg \min_{\alpha} f(\alpha) + g(\alpha)$ with

$$\arg \min_{\alpha} \frac{L}{2} \left\| \alpha - \left(\alpha^{k-1} - \frac{1}{L} \nabla f(\alpha^{k-1}) \right) \right\|^2 + g(\alpha)$$

ISTA

To optimize $F(\boldsymbol{\alpha}) = f(\boldsymbol{\alpha}) + g(\boldsymbol{\alpha})$, we adopt the following approximation model, and define

$$Q_L(\boldsymbol{\alpha}, \boldsymbol{\alpha}^{k-1}) = f(\boldsymbol{\alpha}^{k-1}) + (\boldsymbol{\alpha} - \boldsymbol{\alpha}^{k-1})^T \nabla f(\boldsymbol{\alpha}^{k-1}) + \frac{L}{2} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}^{k-1}\|^2 + g(\boldsymbol{\alpha})$$

which admits a unique minimizer,

$$P_L(\boldsymbol{\alpha}^{k-1}) = \arg \min_{\boldsymbol{\alpha}} \{Q_L(\boldsymbol{\alpha}, \boldsymbol{\alpha}^{k-1})\}$$

The basic step of ISTA for the above problem thus reduces to

$$\boldsymbol{\alpha}^k = P_L(\boldsymbol{\alpha}^{k-1})$$

ISTA

The objective function of ISTA has the form of

$$\arg \min F(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 = f(\boldsymbol{\alpha}) + g(\boldsymbol{\alpha})$$

ISTA with constant stepsize

Input: L : Lipschitz constant of ∇f .

Step 0: Take $\boldsymbol{\alpha}^0 \in \mathbb{R}^n$

Step k: ($k \geq 1$) Compute

$$\boldsymbol{\alpha}^k = P_L(\boldsymbol{\alpha}^{k-1})$$

ISTA

Recall the objective function of ISTA,

$$\arg \min F(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 = f(\boldsymbol{\alpha}) + \lambda g(\boldsymbol{\alpha})$$

For the function $f(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\alpha} - \mathbf{y}\|^2$, we have $\nabla f(\boldsymbol{\alpha}) = \mathbf{X}^T (\mathbf{X}\boldsymbol{\alpha} - \mathbf{y})$

Then we have

$$Q_L(\boldsymbol{\alpha}, \boldsymbol{\alpha}^{k-1}) = \frac{L}{2} \left\| \boldsymbol{\alpha} - \left(\boldsymbol{\alpha}^{k-1} - \frac{1}{L} \mathbf{X}^T (\mathbf{X}\boldsymbol{\alpha}^{k-1} - \mathbf{y}) \right) \right\|^2 + \lambda \|\boldsymbol{\alpha}\|_1$$

As a result, we have the equivalent problem,

$$\boldsymbol{\alpha}^k = P_L(\boldsymbol{\alpha}^{k-1}) = \arg \min \frac{L}{2} \|\boldsymbol{\alpha} - \theta(\boldsymbol{\alpha}^{k-1})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$$

$$\theta(\boldsymbol{\alpha}^{k-1}) = \boldsymbol{\alpha}^{k-1} - \frac{1}{L} \mathbf{X}^T (\mathbf{X}\boldsymbol{\alpha}^{k-1} - \mathbf{y})$$

ISTA

$$\begin{aligned}\boldsymbol{\alpha}^k &= P_L(\boldsymbol{\alpha}^{k-1}) = \operatorname{argmin} \frac{L}{2} \|\boldsymbol{\alpha} - \theta(\boldsymbol{\alpha}^{k-1})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \\ \theta(\boldsymbol{\alpha}^{k-1}) &= \boldsymbol{\alpha}^{k-1} - \frac{1}{L} \mathbf{X}^T (\mathbf{X} \boldsymbol{\alpha}^{k-1} - \mathbf{y})\end{aligned}$$

Remember that for the problem with L1 norm,

$$\boldsymbol{\alpha}^* = \min_{\boldsymbol{\alpha}} h(\boldsymbol{\alpha}) = \lambda \|\boldsymbol{\alpha}\|_1 + \frac{1}{2} \|\boldsymbol{\alpha} - \mathbf{s}\|^2$$

We have the solution as $\boldsymbol{\alpha}^* = \textit{shrink}(\mathbf{s}, \lambda)$, where $\textit{shrink}(\mathbf{s}, \lambda)_j = \textit{sign}(\mathbf{s}_j) \max\{|\mathbf{s}_j| - \lambda, 0\}$.

Therefore, the solution of $\boldsymbol{\alpha}^k = \operatorname{argmin} \frac{L}{2} \|\boldsymbol{\alpha} - \theta(\boldsymbol{\alpha}^{k-1})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$

Is?

ISTA

$$\begin{aligned}\boldsymbol{\alpha}^k &= P_L(\boldsymbol{\alpha}^{k-1}) = \operatorname{argmin} \frac{L}{2} \|\boldsymbol{\alpha} - \theta(\boldsymbol{\alpha}^{k-1})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \\ \theta(\boldsymbol{\alpha}^{k-1}) &= \boldsymbol{\alpha}^{k-1} - \frac{1}{L} \mathbf{X}^T (\mathbf{X} \boldsymbol{\alpha}^{k-1} - \mathbf{y})\end{aligned}$$

Remember that for the problem with L1 norm,

$$\boldsymbol{\alpha}^* = \min_{\boldsymbol{\alpha}} h(\boldsymbol{\alpha}) = \lambda \|\boldsymbol{\alpha}\|_1 + \frac{1}{2} \|\boldsymbol{\alpha} - \mathbf{s}\|^2$$

We have the solution as $\boldsymbol{\alpha}^* = \textit{shrink}(\mathbf{s}, \lambda)$, where $\textit{shrink}(\mathbf{s}, \lambda)_j = \textit{sign}(\mathbf{s}_j) \max\{|\mathbf{s}_j| - \lambda, 0\}$.

Therefore, the solution of $\boldsymbol{\alpha}^k = \operatorname{argmin} \frac{L}{2} \|\boldsymbol{\alpha} - \theta(\boldsymbol{\alpha}^{k-1})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$

is $\boldsymbol{\alpha}^k = \textit{shrink}(\theta(\boldsymbol{\alpha}^{k-1}), \lambda/L)$.

ISTA with backtracking

Drawback: the Lipschitz constant L is not always easily computable (for large-scale problems).

ISTA with backtracking

Step 0. Take $L_0 > 0$, some $\eta > 1$, and $\alpha^0 \in \mathbb{R}^n$

Step k. ($k \geq 1$) Find the **smallest nonnegative integers** i_k such that with $\bar{L} = \eta^{i_k} L_{k-1}$

$$F\left(P_{\bar{L}}(\alpha^{k-1})\right) \leq Q_{\bar{L}}\left(P_{\bar{L}}(\alpha^{k-1}), \alpha^{k-1}\right)$$

Set $L_k = \eta^{i_k} L_{k-1}$ and compute

$$\alpha^k = P_{L_k}(\alpha^{k-1})$$

$$P_L(\alpha^{k-1}) = \arg \min_{\alpha} \{Q_L(\alpha, \alpha^{k-1})\} \quad \alpha^k = P_L(\alpha^{k-1})$$

FISTA (Fast ISTA)

Main Difference: The $P_L(\cdot)$ is not employed on the previous point α^{k-1} , but rather at the point \mathbf{y}^k which uses a very specific linear combination of the previous two points $\{\alpha^{k-1}, \alpha^{k-2}\}$.

FISTA with constant stepsize

Input: L : Lipschitz constant of ∇f .

Step 0. Take $\mathbf{y}^1 = \alpha^0 \in \mathbb{R}^n, t_1 = 1$

Step k. ($k \geq 1$) Compute

$$\alpha^k = P_L(\mathbf{y}^k)$$

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$$

$$\mathbf{y}^{k+1} = \alpha^k + \left(\frac{t_k - 1}{t_{k+1}} \right) (\alpha^k - \alpha^{k-1})$$

FISTA

FISTA with backtracking

Step 0. Take $L_0 > 0$, some $\eta > 1$, and $\mathbf{y}^1 = \boldsymbol{\alpha}^0 \in \mathbb{R}^n$, $t_1 = 1$

Step k. ($k \geq 1$) Find the **smallest nonnegative integers** i_k such that with $\bar{L} = \eta^{i_k} L_{k-1}$

$$F\left(P_{\bar{L}}(\mathbf{y}^k)\right) \leq Q_{\bar{L}}(P_{\bar{L}}(\mathbf{y}^k), \mathbf{y}^k)$$

Set $L_k = \eta^{i_k} L_{k-1}$ and compute

$$\begin{aligned}\boldsymbol{\alpha}^k &= P_L(\mathbf{y}^k) \\ t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2} \\ \mathbf{y}^{k+1} &= \boldsymbol{\alpha}^k + \left(\frac{t_k - 1}{t_{k+1}}\right) (\boldsymbol{\alpha}^k - \boldsymbol{\alpha}^{k-1})\end{aligned}$$

Conclusions

Potential benefits

- To improve accuracy
- Reduce computation
- Reduce space
- Reduce cost of future measurements
- Improved data/model understanding

Conclusions

- **Wrappers methods**

- Use machine learning algorithm as **black box** to find best subset of features.
- Generally infeasible on the model 'big data' problem.

- **Filters methods**

- Features selected **before** machine learning algorithm is run.

- **Embedded methods**

- Feature selection occurs naturally as **part of** the machine learning algorithm.

Conclusions

- **Overfitting and Underfitting**
- **L1-norm and L2-norm**
- **Sparse Coding**
- **Shrink Operator**
- **ISTA and FISTA (optional)**