

Non-linear hypotheses. Overfitting and regularization

Sriram Sankararaman

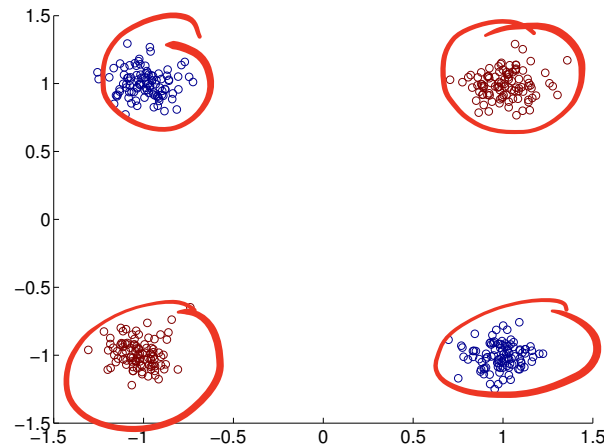
The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Outline

- 1 Nonlinear hypotheses
- 2 Basic ideas to overcome overfitting
- 3 Regularized linear regression (ridge regression)
- 4 A general view of supervised learning

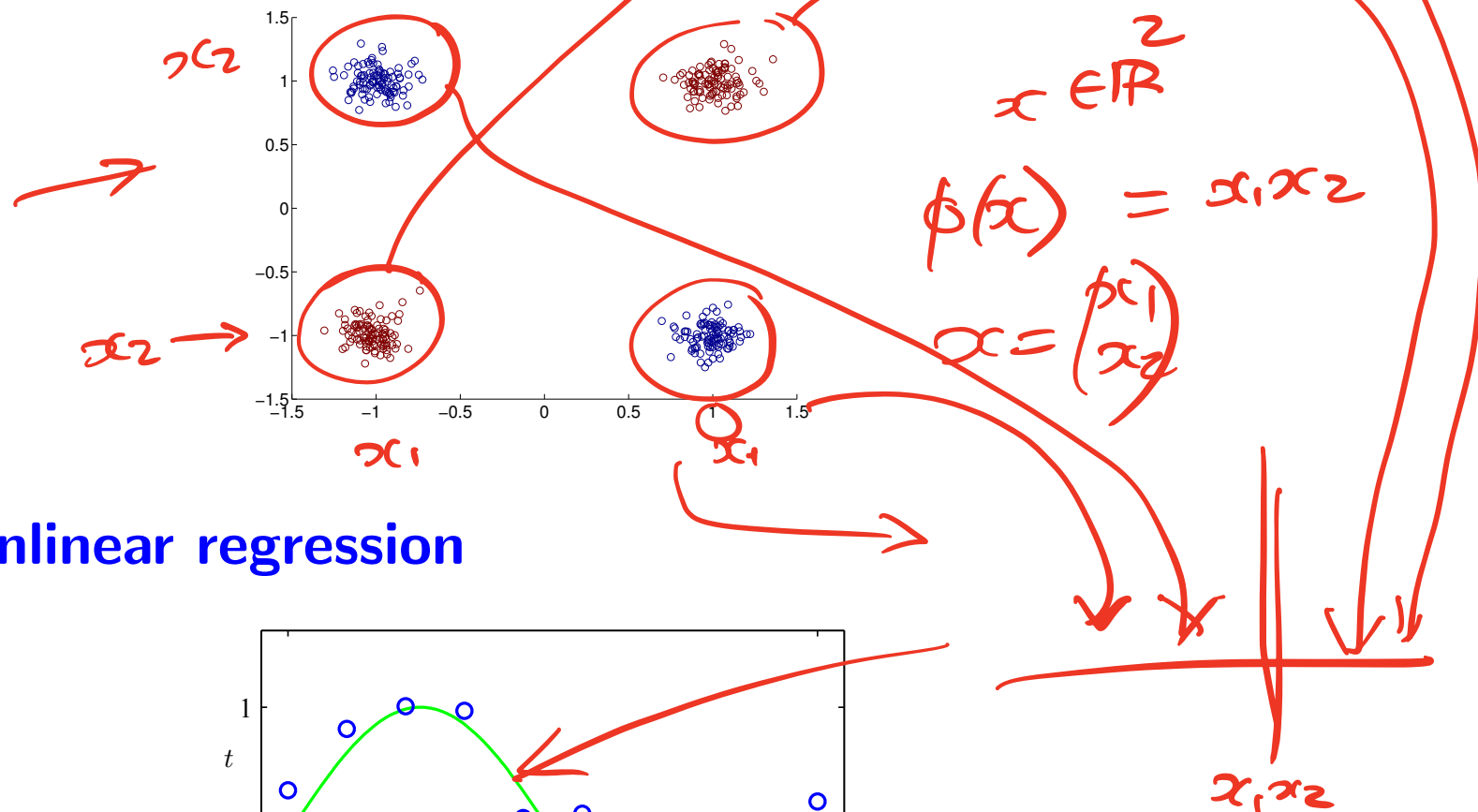
What if data is not linearly separable or fits to a line

Example of nonlinear classification

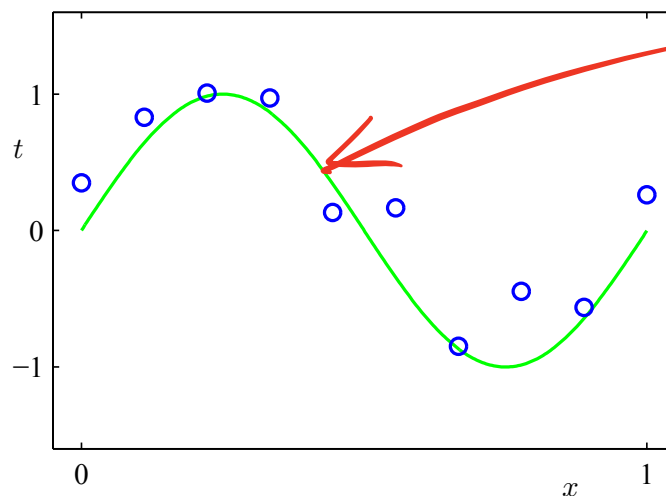


What if data is not linearly separable or fits to a line

Example of nonlinear classification



Example of nonlinear regression



Nonlinear basis for classification

Transform the input/feature

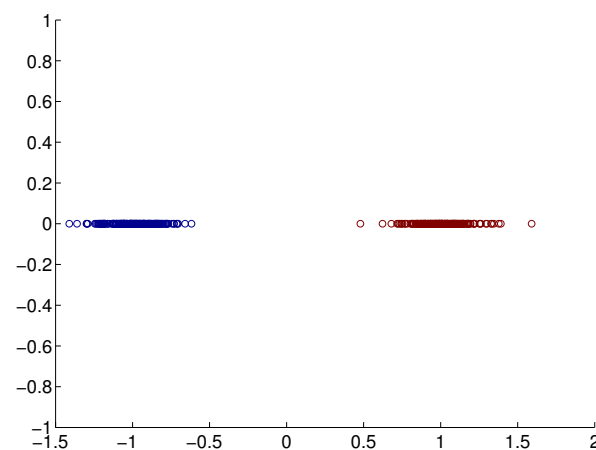
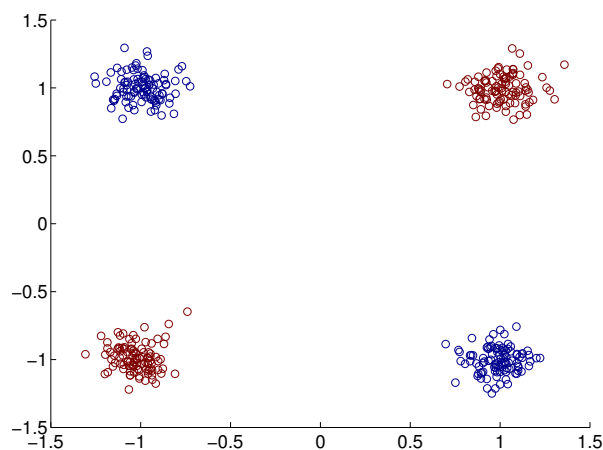
$$\underline{\underline{\phi(\boldsymbol{x})}} : \boldsymbol{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Nonlinear basis for classification

Transform the input/feature

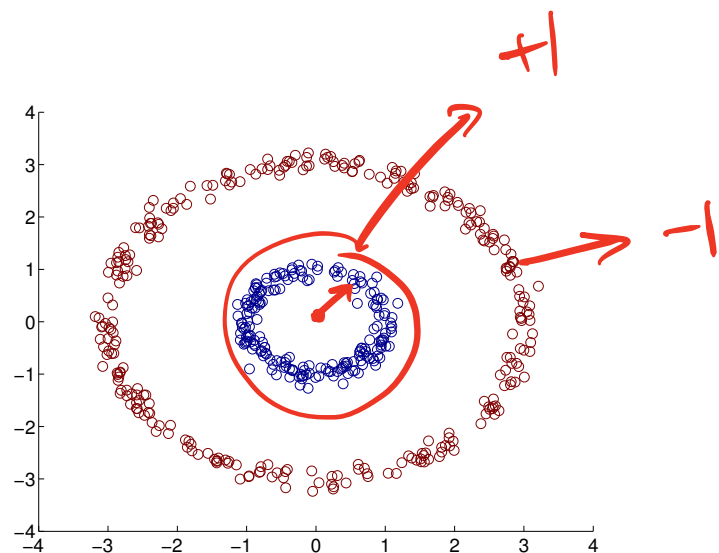
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Transformed training data: linearly separable!



Another example

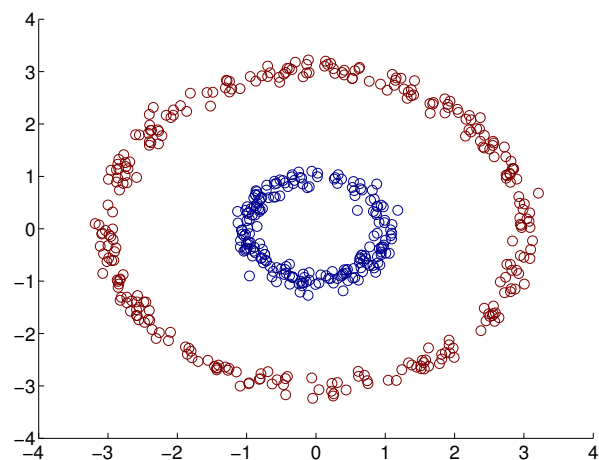
How to transform the input/feature?



$\phi(x)$?

Another example

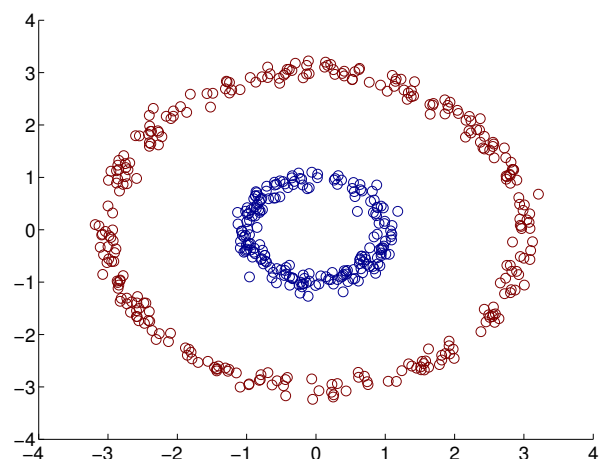
How to transform the input/feature?



$$\phi(\underline{x}) : \underline{x} \in \mathbb{R}^2 \rightarrow \underline{z} = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

Another example

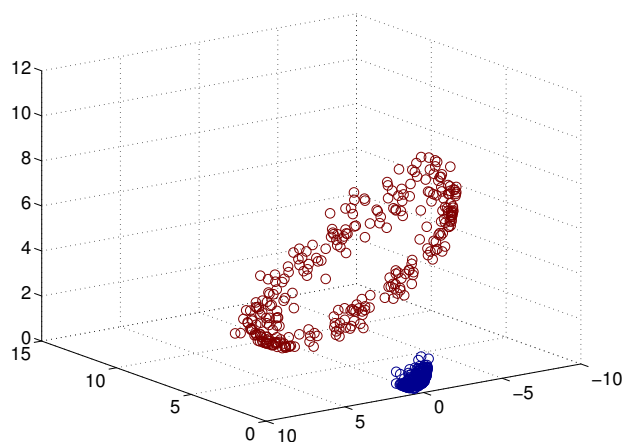
How to transform the input/feature?



$$\phi(x) : x \in \mathbb{R}^2 \rightarrow z = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

$$\phi(x) = \|x\|_2^2$$

Transformed training data: linearly separable



General nonlinear basis functions

We can use a nonlinear mapping

$$\underline{\phi(\underline{x})} : \underline{x} \in \mathbb{R}^D \rightarrow \underline{z} \in \mathbb{R}^M$$

where M is the dimensionality of the new feature/input \underline{z} (or $\phi(\underline{x})$). Note that M could be either greater than D or less than or the same.

General nonlinear basis functions

We can use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

where M is the dimensionality of the new feature/input \mathbf{z} (or $\phi(\mathbf{x})$). Note that M could be either greater than D or less than or the same.

With the new features, we can apply our learning techniques to minimize our errors on the transformed training data

- linear methods: prediction is based on $\theta^T \phi(\mathbf{x})$
- other methods: nearest neighbors, decision trees, etc

Regression with nonlinear basis

Residual sum squares

$$\underset{\theta}{\operatorname{argmin}} J(\theta) = \sum_n [\theta^T \phi(x_n) - y_n]^2$$

Handwritten notes: $\theta^T x_n$ with an arrow pointing to $\phi(x_n)$; $\underset{w}{\operatorname{argmin}} J(w) = \sum_n (w^T x_n - y_n)^2$

where $\theta \in \mathbb{R}^M$, the same dimensionality as the transformed features $\phi(x)$.

Regression with nonlinear basis

Residual sum squares

$$\sum_n [\theta^T \underline{\phi(x_n)} - y_n]^2$$

$$\underset{\omega}{\operatorname{argmin}} J(\omega) = \sum_n (\omega^T x_n - y_n)^2$$
$$\hat{\omega} = (\underline{X^T X})^{-1} \underline{X^T y}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

where $\theta \in \mathbb{R}^M$, the same dimensionality as the transformed features $\phi(x)$.

The linear regression solution can be formulated with the new design matrix

$$M \ll D$$

$$\underline{\Phi} = \begin{pmatrix} \Rightarrow \phi(x_1)^T \\ \Rightarrow \phi(x_2)^T \\ \vdots \\ \Rightarrow \phi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}, \quad \underline{\hat{\theta}} = (\underline{\Phi^T \Phi})^{-1} \underline{\Phi^T y}$$

$\underline{X} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}$

Example with regression

Polynomial basis functions

$x \in \mathbb{R}$
 $y \in \mathbb{R}$

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow \underline{h(x)} = \underline{\theta_0} + \sum_{m=1}^M \underline{\theta_m} x^m$$

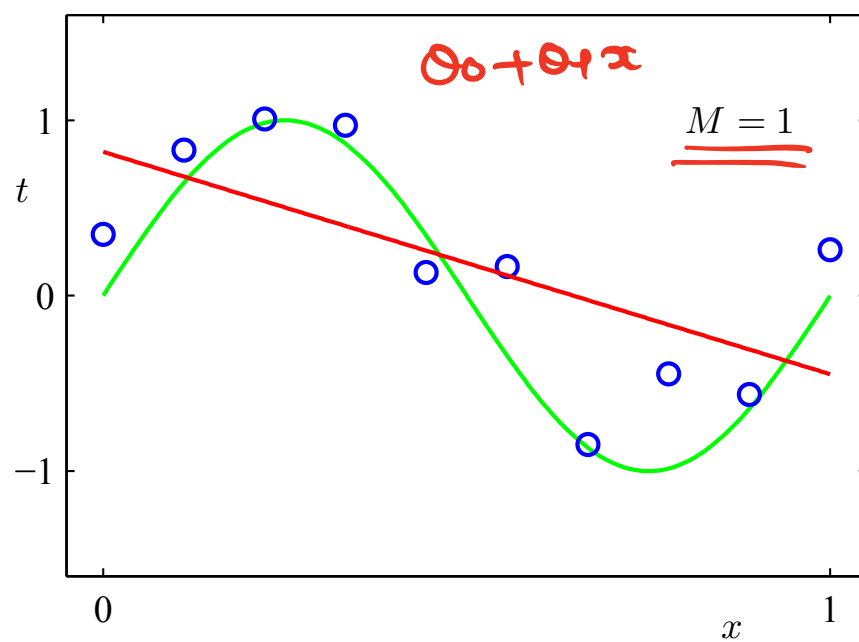
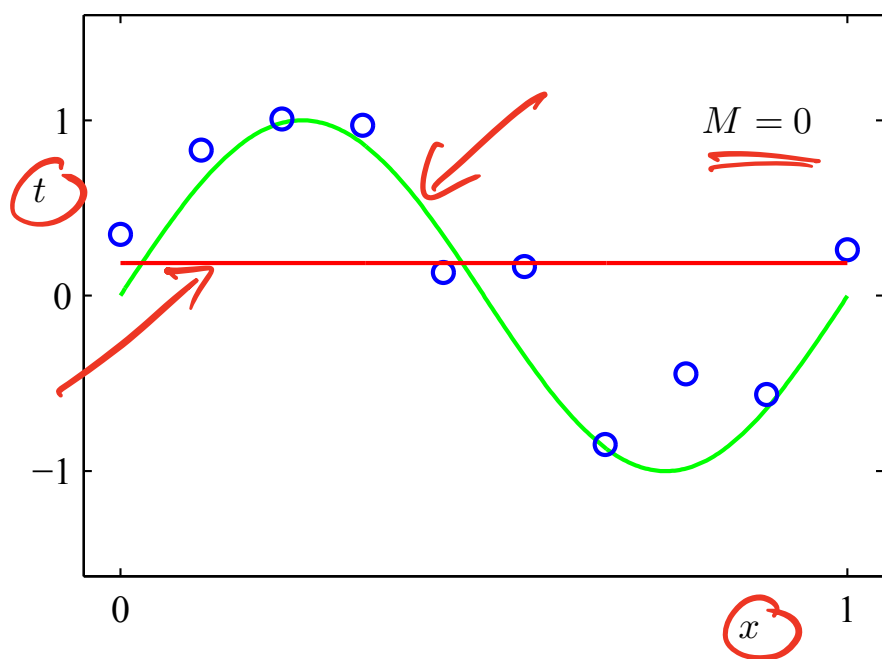
$(M+1)$

Example with regression

Polynomial basis functions

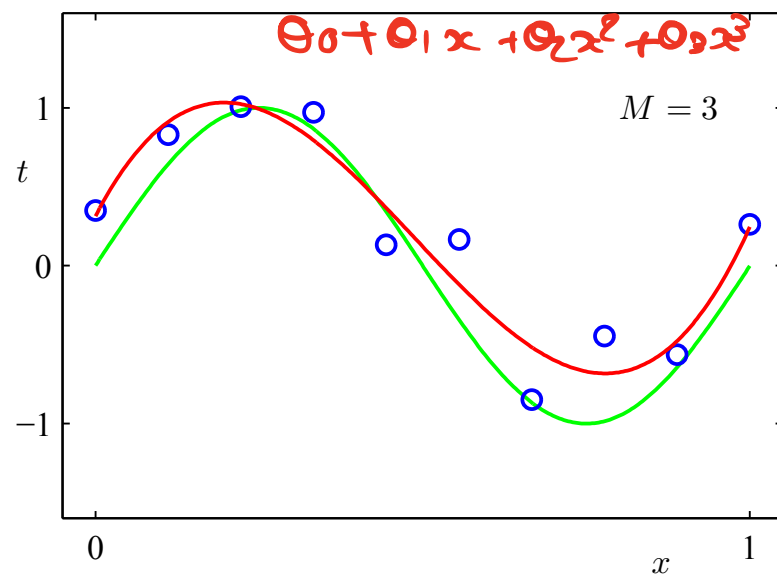
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow h(x) = \theta_0 + \sum_{m=1}^M \theta_m x^m$$

Fitting samples from a sine function: *underrfitting* as $h(x)$ is too simple



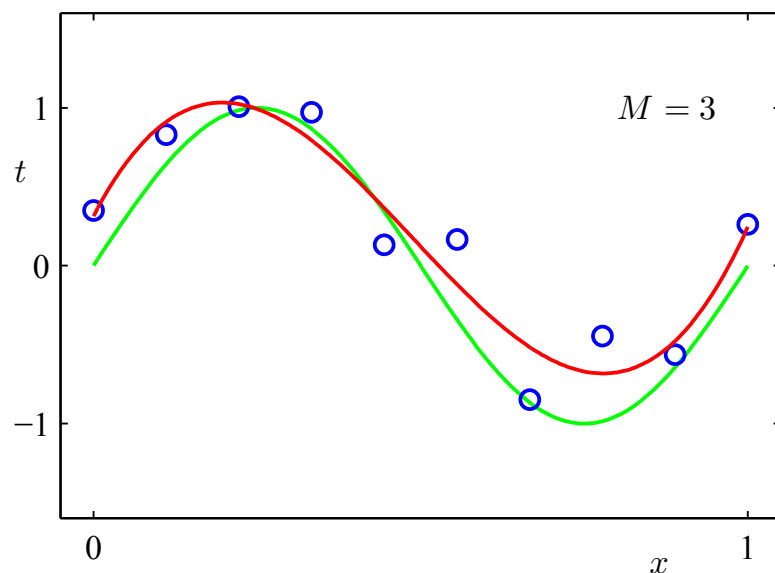
Adding high-order terms

M=3

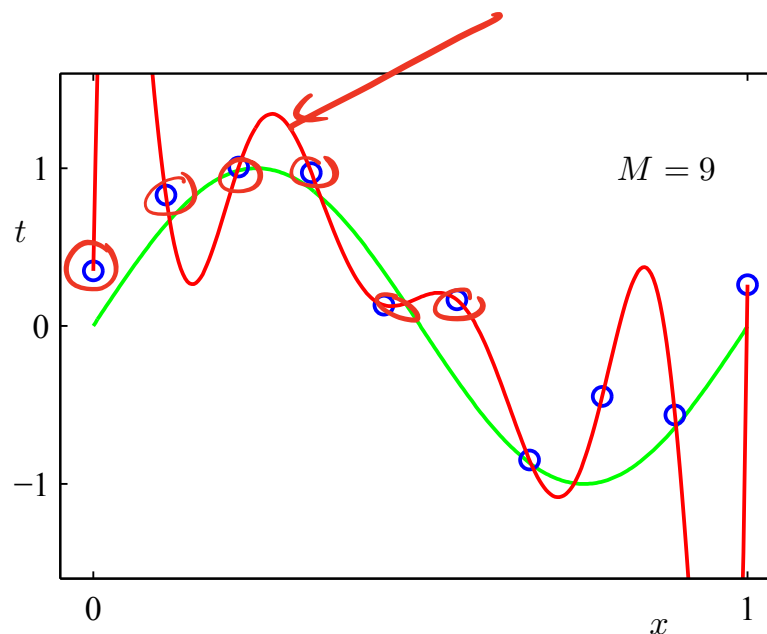


Adding high-order terms

M=3



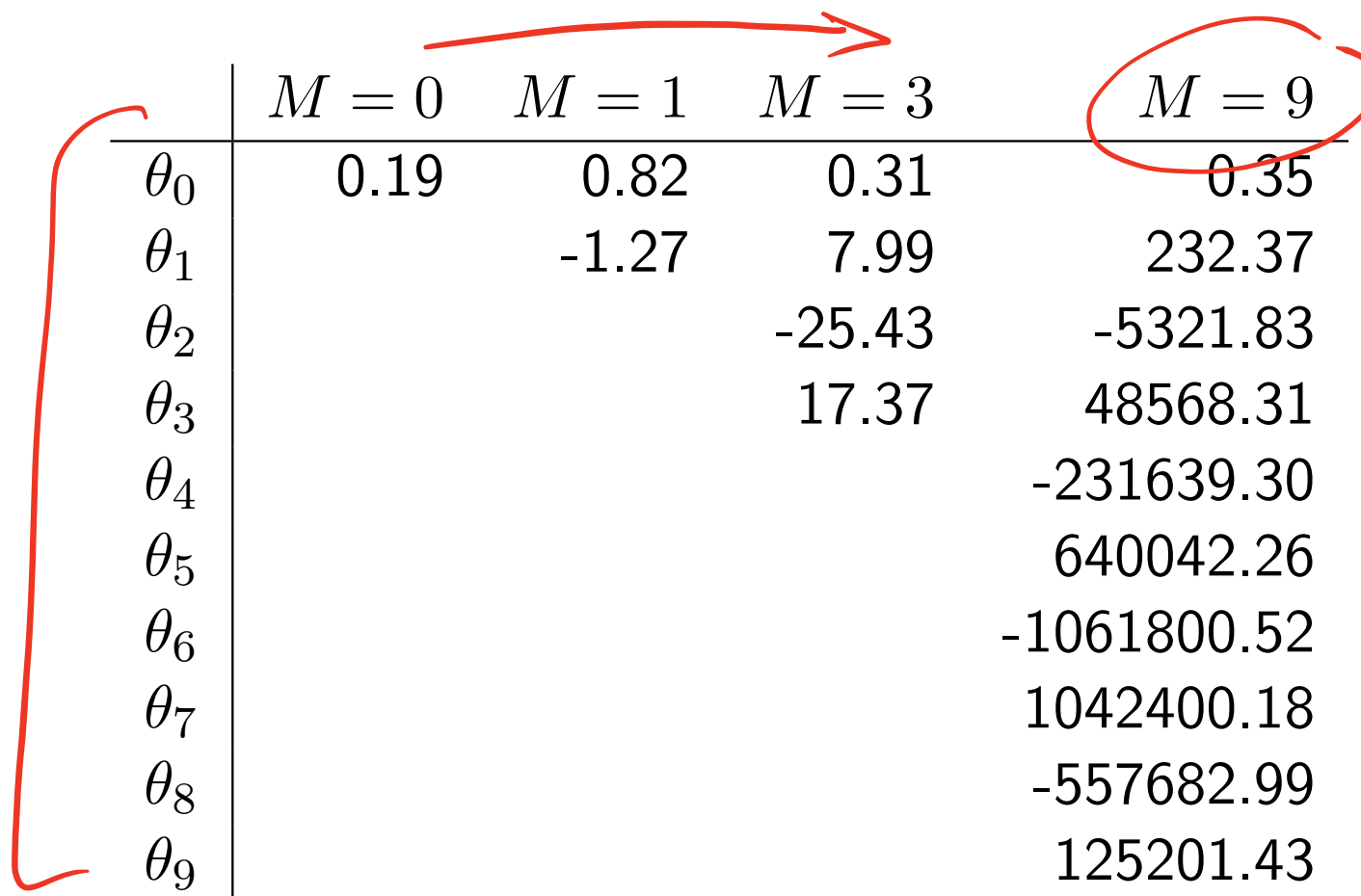
M=9: *overfitting*



More complex features lead to better results on the training data, but potentially worse results on new data, e.g., test data!

Overfitting

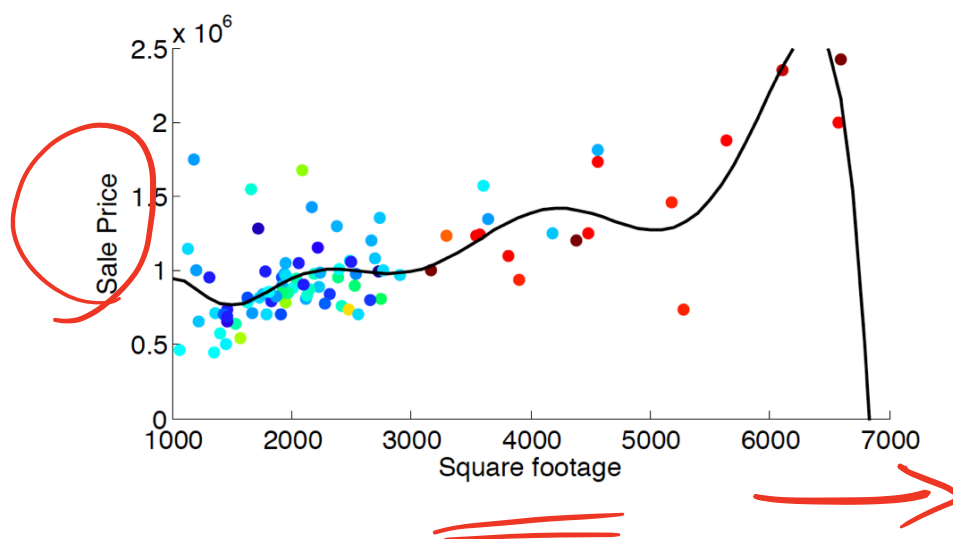
Parameters for higher-order polynomials are very large



	$M = 0$	$M = 1$	$M = 3$	$M = 9$
θ_0	0.19	0.82	0.31	0.35
θ_1		-1.27	7.99	232.37
θ_2			-25.43	-5321.83
θ_3			17.37	48568.31
θ_4				-231639.30
θ_5				640042.26
θ_6				-1061800.52
θ_7				1042400.18
θ_8				-557682.99
θ_9				125201.43

Overfitting can be quite disastrous

Fitting the housing price data with $M = 3$

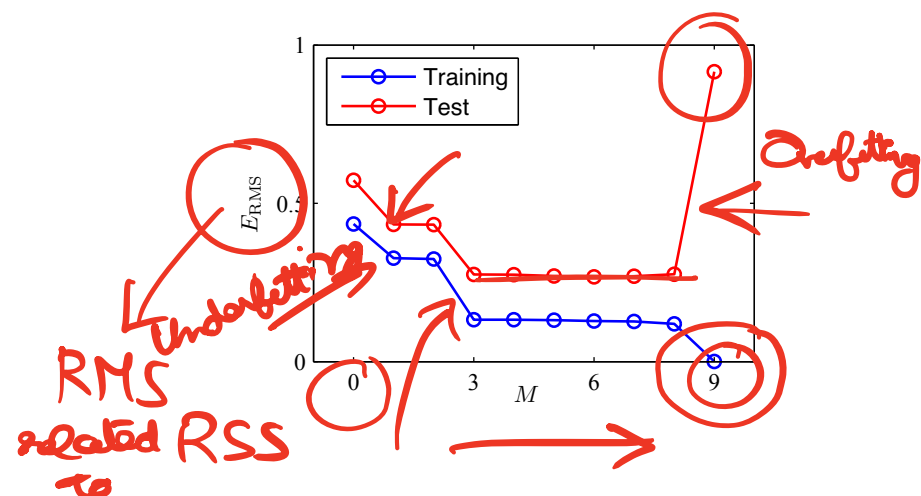


Note that the price would go to zero (or negative) if you buy bigger ones!
This is called poor generalization/overfitting.

Detecting overfitting

Plot model complexity versus objective function

As a model increases in complexity, performance on training data keeps improving while performance on test data may first improve but eventually deteriorate.

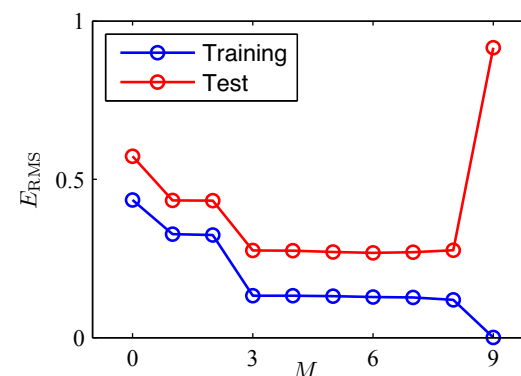


- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.

Detecting overfitting

Plot model complexity versus objective function

As a model increases in complexity, performance on training data keeps improving while performance on test data may first improve but eventually deteriorate.



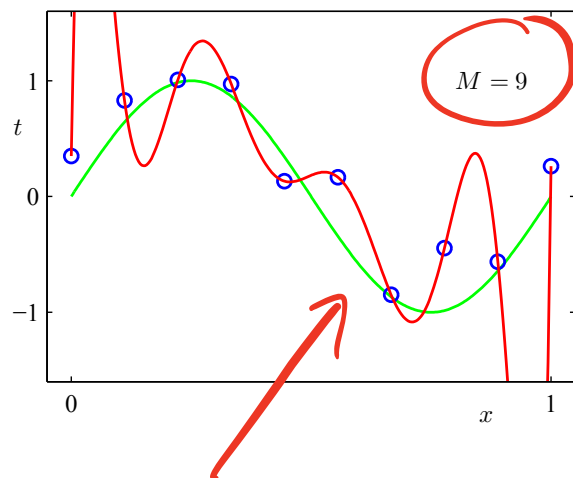
- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.
- Vertical axis:
 - 1 For regression, residual sum of squares or residual mean squared (squared root of RSS)
 - 2 For classification, classification error rate.

Outline

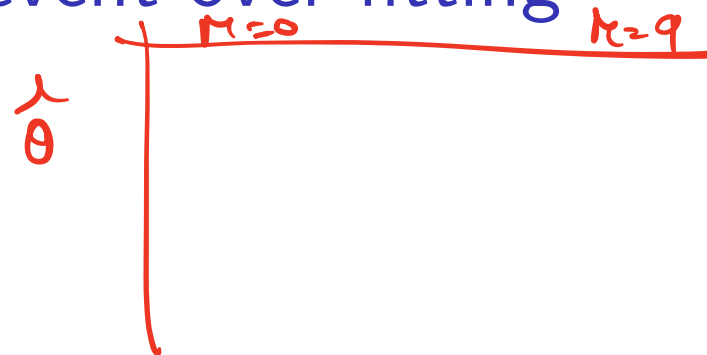
- 1 Nonlinear hypotheses
- 2 Basic ideas to overcome overfitting
 - Use more training data
 - Regularization methods
 - Regularized classification
- 3 Regularized linear regression (ridge regression)
- 4 A general view of supervised learning

Use more training data to prevent over fitting

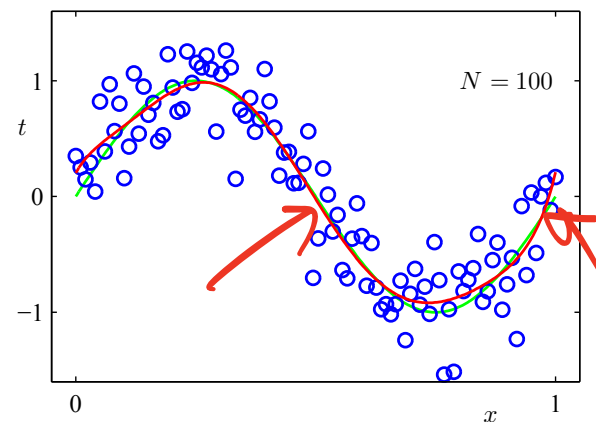
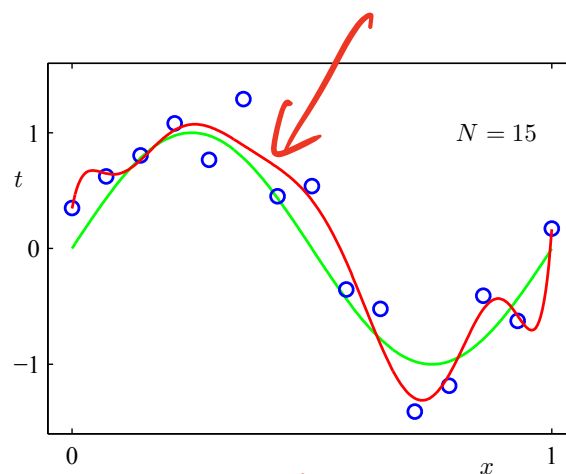
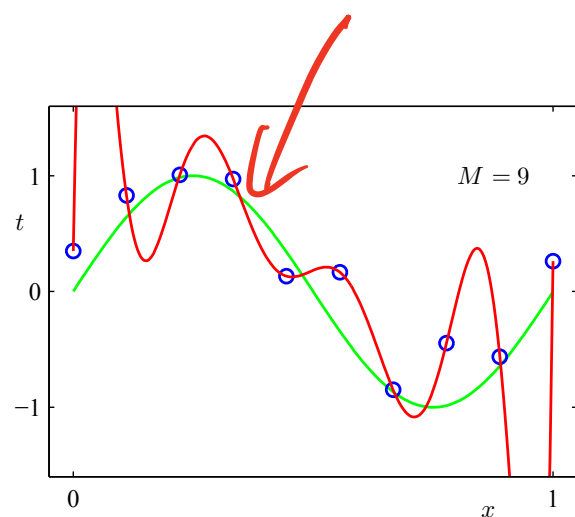
The more, the merrier



Use more training data to prevent over fitting

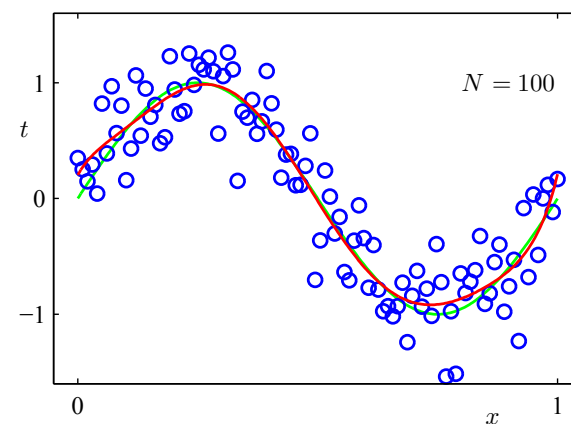
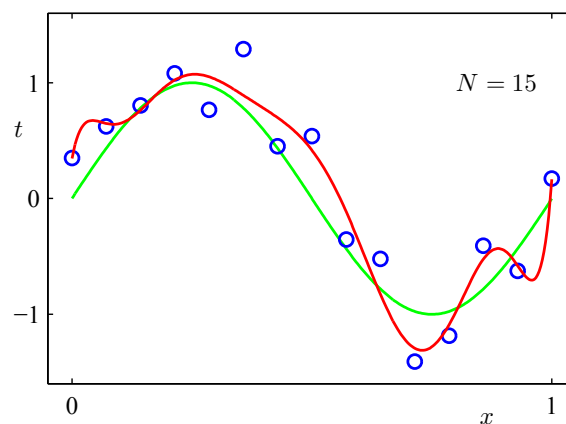
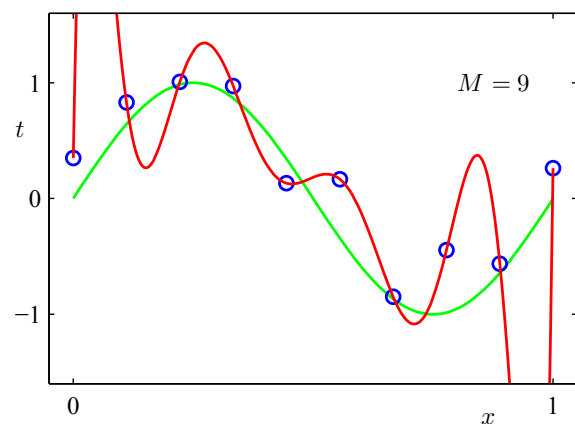


The more, the merrier



Use more training data to prevent over fitting

The more, the merrier



What if we do not have a lot of data?

Regularization methods

Intuition: For a linear model for regression

$$\underline{w^T x + b}$$

$$w = (w_1, 0, w_3, \dots, 0)$$

we can try to identify 'simpler' models. But what does it mean for a model to be *simple*?

- 1) Fewer parameters
10 ← 100
Simpler than a
- 2) Parameters take smaller values
 $|w| < 1$
 $|w| > 10^6$
- 3) Requires less computation
- 4) Most features are ignored
= most entries $w = 0$
- 5) Smoother
- 6) Easier to interpret

Regularization methods

Intuition: For a linear model for regression

$$\underline{w^T x + b}$$

we can try to identify 'simpler' models. But what does it mean for a model to *be simple*?

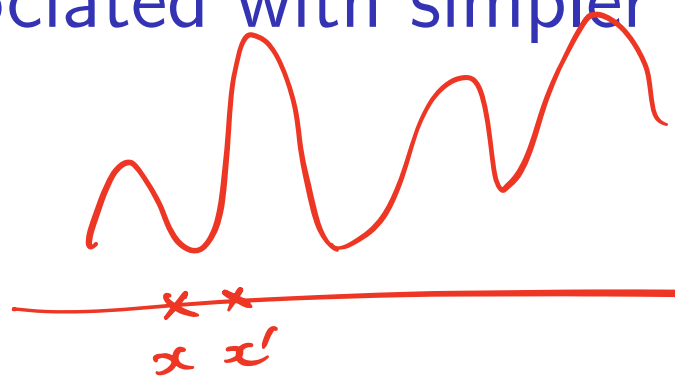
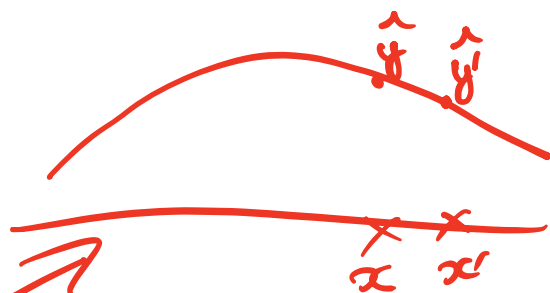
Assumption (inductive bias)

A simpler model is one where most of the weights are zero.

A simpler model is one with smaller weights.


Smoothness

Why are smaller weights associated with simpler models?



Simpler functions are smoother, *i.e.*, nearby values of x have similar outputs \hat{y} .

Two values x and x' that differ in the first component by a small value ϵ .

Their predictions \hat{y} and \hat{y}' differ by ϵw_1 .

Smaller w_1 (closer to zero), more similar are the predictions.

$$\begin{aligned}\hat{y} &= w_1 x_1 + \cancel{w_2 x_2} + \dots + \cancel{w_D x_D} \\ \hat{y}' &= w_1 x_1' + \cancel{w_2 x_2} + \dots + \cancel{w_D x_D} \\ \hat{y} - \hat{y}' &= w_1 (x_1 - x_1') \\ &= \boxed{w_1 \epsilon}\end{aligned}$$

Regularized linear regression

A new cost function or error function to minimize

$$J(\underline{\mathbf{w}}, b) = \sum_n (y_n - \underline{\mathbf{w}}^T \mathbf{x}_n - b)^2 + \lambda \|\underline{\mathbf{w}}\|_2^2$$

Handwritten notes: $\lambda > 0$, $(w_1^2 + w_2^2 + \dots + w_n^2)$, $\lambda \|\mathbf{w}\|_2^2$ is circled and labeled with an arrow pointing to the text "regularization/regularizer". The term \sum_n is underlined and labeled "RSS".

where $\lambda > 0$. This extra term $\|\underline{\mathbf{w}}\|_2^2$ is called regularization/regularizer and controls the model complexity.

$$\lambda \uparrow \infty$$

$$\hat{\mathbf{w}}$$

$$\lambda \|\mathbf{w}\|_2^2$$

$$\hat{\mathbf{w}} \Rightarrow 0$$

$$\lambda \downarrow 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_{OLS}$$

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then

$$\hat{\mathbf{w}} \rightarrow \mathbf{0}$$

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then

$$\hat{\mathbf{w}} \rightarrow \mathbf{0}$$

- If $\lambda \rightarrow 0$, then we trust our data more. Numerically,

$$\hat{\mathbf{w}} \rightarrow \arg \min \sum_n (\mathbf{w}^T \mathbf{x}_n + b - y_n)^2$$

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the OLS (Ordinary Least Squares) solution

$$\arg \min \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2 \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and reduces to the OLS solution when $\lambda = 0$, as expected.

$$\nabla(\text{OLS}) + \nabla(\quad) = 0$$

$$\begin{matrix} \times \\ N \times D \end{matrix}$$

$$\mathbf{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{pmatrix} \begin{matrix} D \times D \\ D \times D \\ D \times D \end{matrix} + \lambda \begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix} \begin{matrix} D \times D \\ D \times D \\ D \times D \end{matrix}$$

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the OLS (Ordinary Least Squares) solution

$$\arg \min \sum_n (y_n - \underline{\mathbf{w}^T \mathbf{x}_n - b})^2 + \lambda \underline{\|\mathbf{w}\|_2^2} \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and reduces to the OLS solution when $\lambda = 0$, as expected.

If we have to use numerical procedure, the gradient would change nominally too,

$$\underline{\nabla J(\mathbf{w})} = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{w})$$

As long as $\lambda \geq 0$, the optimization is convex.

RSS

Example: fitting data with polynomials

Our regression model

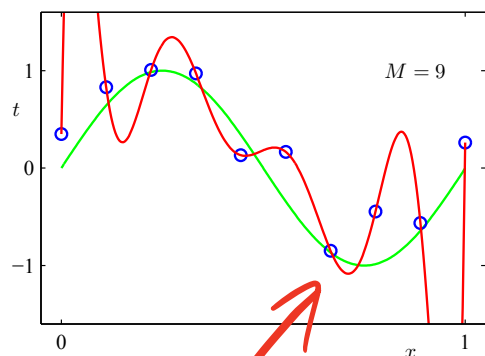
$$y = \sum_{m=1}^M w_m x^m$$

Regularization would discourage large parameter values as we saw with the OLS solution, thus potentially preventing overfitting.

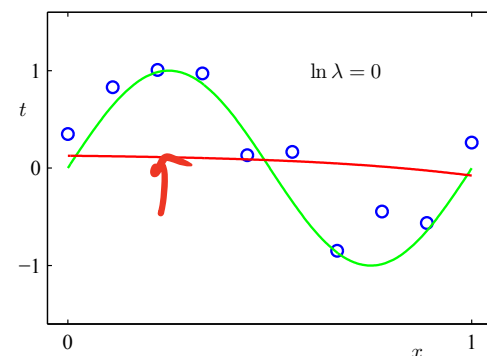
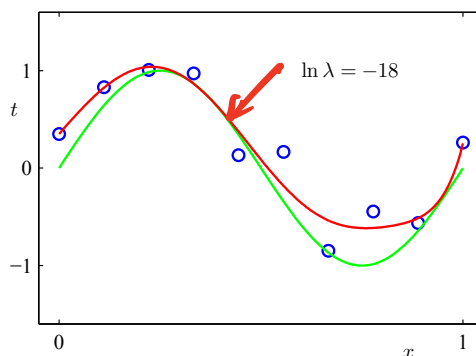
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0	0.19	0.82	0.31	0.35
w_1		-1.27	7.99	232.37
w_2			-25.43	-5321.83
w_3			17.37	48568.31
w_4				-231639.30
w_5				640042.26
w_6				-1061800.52
w_7				1042400.18
w_8				-557682.99
w_9				125201.43

Overfitting in terms of λ

Overfitting is reduced from complex model to simpler one with the help of increasing regularizer

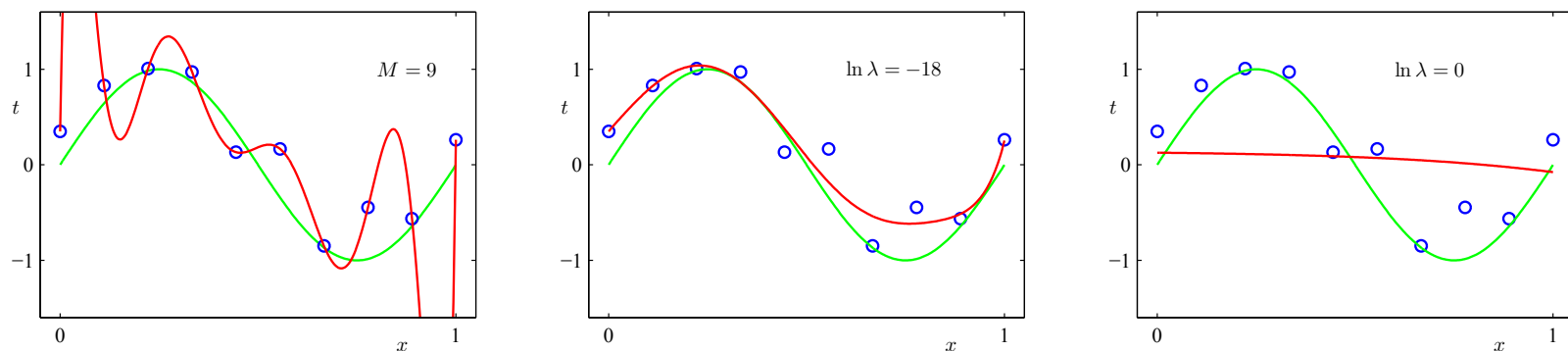


$\lambda=0$

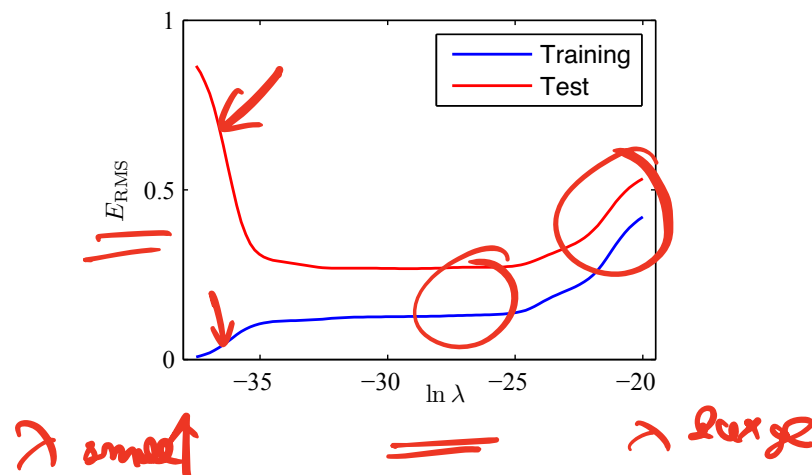


Overfitting in terms of λ

Overfitting is reduced from complex model to simpler one with the help of increasing regularizer



λ vs. residual error shows the difference of the model performance on training and testing dataset



The effect of λ

Large λ attenuates parameters towards 0

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0	0.35	0.35	0.13
w_1	232.37	4.74	-0.05
w_2	-5321.83	-0.77	-0.06
w_3	48568.31	-31.97	-0.06
w_4	-231639.30	-3.89	-0.03
w_5	640042.26	55.28	-0.02
w_6	-1061800.52	41.32	-0.01
w_7	1042400.18	-45.95	-0.00
w_8	-557682.99	-91.53	0.00
w_9	125201.43	72.68	0.01

l_2 regularized logistic regression

Adding regularizer to the cost function for logistic regression

$$J(\mathbf{w}, b) = - \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$
$$h_{\mathbf{w}, b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

l_2 regularized logistic regression

Adding regularizer to the cost function for logistic regression

$$J(\mathbf{w}, b) = - \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

$$h_{\mathbf{w}, b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

Numerical optimization

- Objective functions remains convex as long as $\lambda \geq 0$.
- Gradients and Hessians are changed marginally and can be easily derived.

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

$$\min_{\lambda} \min_{(w,b)} \left[\text{RSS}(w,b) + \lambda \|w\|_2^2 \right]$$

$$\lambda = 0$$

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

No: as this will set λ to zero, i.e., without regularization, defeating our intention to use it to control model complexity and to gain better generalization.

λ is thus a hyperparameter. To tune it,

- We can use a development/validation dataset independent of training and testing dataset.

The procedure is similar to choosing K in the nearest neighbor classifiers.

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

No: as this will set λ to zero, i.e., without regularization, defeating our intention to use it to control model complexity and to gain better generalization.

λ is thus a **hyperparameter**. To tune it,

- We can use a development/validation dataset independent of training and testing dataset.

The procedure is similar to choosing K in the nearest neighbor classifiers.

For different λ , we get $\hat{\mathbf{w}}$ and evaluate the model on the development/validation dataset.

We then plot the curve λ versus prediction error (accuracy, classification error) and find the place that the performance on the validation is the best.

Outline

- 1 Nonlinear hypotheses
- 2 Basic ideas to overcome overfitting
- 3 Regularized linear regression (ridge regression)
- 4 A general view of supervised learning

Solution to linear regression

Compute the gradient of J and find the value of θ at which the gradient vanishes.

$$\Rightarrow \nabla J(\theta) = 0$$

This leads us to solve the normal equations

$$Ax=b$$

$$\Rightarrow X^T X \theta = X^T y$$

to obtain

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

X
Design
matrix
($N \times D$)

y
($N \times 1$)
Target
vector

θ
Parameter
($D \times 1$)

What if $X^T X$ is not invertible

Answer 1: $N < D$. Intuitively, not enough data to estimate all the parameters.

Answer 2: X columns are not linearly independent. Intuitively, there are two features that are perfectly correlated. In this case, solution is not unique.

$$\hat{\theta} = \underbrace{(X^T X)^{-1}}_{\text{not invertible}} X^T y \quad \times$$
$$\hat{\theta} = \underbrace{(X^T X + \lambda \mathbb{I})^{-1}}_{\lambda > 0} X^T y$$

l_2 regularized linear regression (ridge regression)

Solution

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T y$$

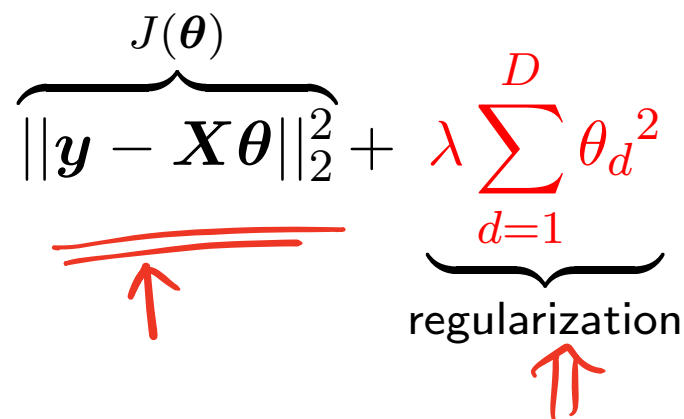
$\underset{\theta}{\operatorname{argmin}} \quad J(\theta) = \underset{\substack{\downarrow \\ \text{OLS cost fn}}}{\text{RSS}(\theta)} + \lambda \|\theta\|_2^2$

l_2 regularized linear regression (ridge regression)

Solution

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

This is equivalent to adding an extra term to $J(\boldsymbol{\theta})$

$$\overbrace{\| \mathbf{y} - \mathbf{X} \boldsymbol{\theta} \|_2^2}^{J(\boldsymbol{\theta})} + \underbrace{\lambda \sum_{d=1}^D \theta_d^2}_{\text{regularization}}$$


- $\sum_{d=1}^D \theta_d^2 = \sum_{d=1}^D w_d^2 = \|\mathbf{w}\|_2^2$
- l_2 -regularization uses the squared 2-norm

l_2 regularized linear regression (ridge regression)

Solution

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

This is equivalent to adding an extra term to $J(\boldsymbol{\theta})$

$$\overbrace{\underbrace{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}_{\text{error}}}^{J(\boldsymbol{\theta})} + \underbrace{\lambda \sum_{d=1}^D \theta_d^2}_{\text{regularization}}$$

Trade-off two quantities: minimize the error while keeping the weights small.

- $\sum_{d=1}^D \theta_d^2 = \sum_{d=1}^D w_d^2 = \|\mathbf{w}\|_2^2$
- l_2 -regularization uses the squared 2-norm

l_2 regularized linear regression (ridge regression)

Solution

$$\hat{\theta} = (\underbrace{X^T X + \lambda I}_{\text{regularization}})^{-1} X^T y$$

This is equivalent to adding an extra term to $J(\theta)$

$$\underbrace{J(\theta)}_{||y - X\theta||_2^2} + \underbrace{\lambda \sum_{d=1}^D \theta_d^2}_{\text{regularization}}$$

$\log(\lambda)$

- $\sum_{d=1}^D \theta_d^2 = \sum_{d=1}^D w_d^2 = \|w\|_2^2$
- l_2 -regularization uses the squared 2-norm

l_1 -regularization
 $\| \theta \|_1 = \sum_{d=1}^D |\theta_d|$

Benefits

- Numerically more stable, invertible matrix
- Prevent overfitting

How to choose λ ?

λ is referred as *hyperparameter*

- In contrast θ is the parameter vector
- Use validation or cross-validation to find good choice of λ

Mini-summary

- l_2 regularized linear regression (ridge regression) can be helpful to avoid numerical instability.
- Only a minor modification to the OLS estimate
- New hyperparameter λ needs to be chosen using cross-validation.

Outline

- 1 Nonlinear hypotheses
- 2 Basic ideas to overcome overfitting
- 3 Regularized linear regression (ridge regression)
- 4 A general view of supervised learning

Supervised learning

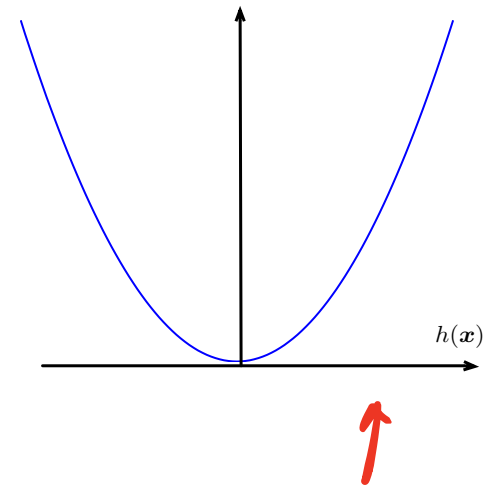
We aim to build a function $\underline{h(\underline{x})}$ to predict the true value y associated with x . If we make a mistake, we incur a *loss*

Loss function $\rightarrow \underline{\ell(y, \underline{h(x)})}$

Example: squared loss function for regression when y is continuous

$$\ell(y, \underline{h(x)}) = \underline{[h(x) - y]^2}$$

Ex: when $\underline{y = 0}$



Other types of loss functions

For classification: 0/1 loss

$$\ell(y, h(\mathbf{x})) = \underbrace{1}_{\text{if } y \neq h(\mathbf{x})} \{y \neq h(\mathbf{x})\}$$

$$\begin{array}{ll} y = h(\mathbf{x}) & \ell(y, h(\mathbf{x})) = 0 \\ y \neq h(\mathbf{x}) & \ell(y, h(\mathbf{x})) = 1 \end{array}$$

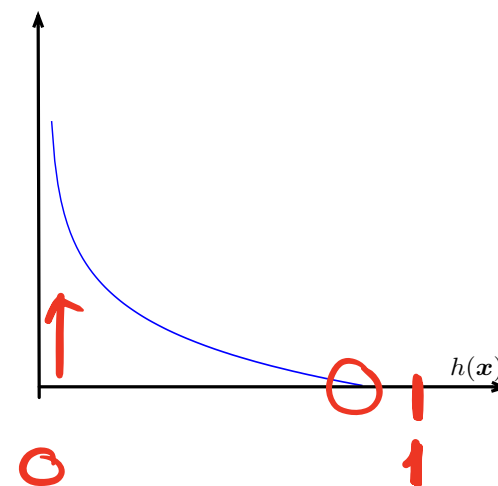
Other types of loss functions

For classification: *logistic* loss

$$\ell(y, h(\mathbf{x})) = -y \log h(\mathbf{x}) - (1-y) \log[1-h(\mathbf{x})]$$

Assumes h outputs values in $[0, 1]$.

Ex: when $y = 1$



Measure how good our hypothesis h is

Risk / expected test loss: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y)$$

$$\begin{aligned} \ell(h_1(\underline{x}), \underline{y}) &< \ell(h_2(\underline{x}), \underline{y}) \\ \ell(h_1(\tilde{x}), \tilde{y}) &> \ell(h_2(\tilde{x}), \tilde{y}) \end{aligned}$$

Measure how good our hypothesis h is

Risk/ expected test loss: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y)$$

However, we cannot compute $\mathcal{R}[h(\mathbf{x})]$, so we use *empirical risk/training error*, given a training dataset \mathcal{D} (\mathbf{x}_n, y_n)

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as $N \rightarrow +\infty$,

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] \rightarrow \mathcal{R}[h(\mathbf{x})]$$

Pick a good hypothesis h

A good hypothesis h is one that minimizes risk/expected test loss $\mathcal{R}[h(\mathbf{x})]$ but we cannot even compute it!

Instead pick h that minimizes empirical risk/training error $\mathcal{R}^{\text{EMP}}[h(\mathbf{x})]$

This strategy is known as **empirical risk minimization**.

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, and we use squared loss
- For logistic regression, $h_{\mathbf{w},b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$, and we use logistic loss
- Finding the best h is achieved by searching for (\mathbf{w}, b) that minimizes the training error/empirical risk.

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, and we use squared loss
- For logistic regression, $h_{\mathbf{w},b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$, and we use logistic loss
- Finding the best h is achieved by searching for (\mathbf{w}, b) that minimizes the training error/empirical risk.

ERM might be problematic

- If $h(\mathbf{x})$ is complicated enough,

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] \rightarrow 0$$

- But then $h(\mathbf{x})$ is unlikely to do well in predicting things out of the training dataset \mathcal{D}
- This is called *poor generalization* or *overfitting*. We have just discussed approaches to address this issue.

Regularizer

Instead of \mathcal{R}^{emp} , use


$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \mathcal{R}^{\text{EMP}}[\mathbf{h}_{\mathbf{w}, b}(\mathbf{x})] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, \mathbf{h}_{\mathbf{w}, b}(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

Loss functions ℓ , $\hat{y} = \mathbf{w}^T \mathbf{x} + b$

- Zero/One: $\ell(y, \hat{y}) = \mathbf{1}\{y \neq \hat{y}\} = \mathbf{1}\{y\hat{y} \leq 0\}$
- Squared: $\ell(y, \hat{y}) = [y - \hat{y}]^2$

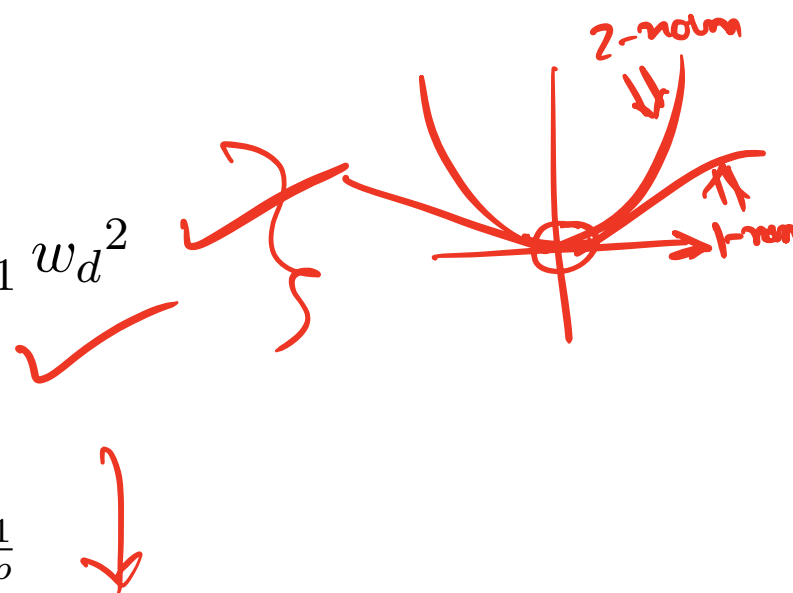
Regularizer

Instead of \mathcal{R}^{emp} , use

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \mathcal{R}^{\text{EMP}}[h_{\mathbf{w}, b}(x)] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, h_{\mathbf{w}, b}(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$


Regularizer

- Squared 2-norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_2^2 = \sum_{d=1}^D w_d^2$
- 1-norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$
- 0-norm: $R(\mathbf{w}, b) = \sum_{d=1}^D \mathbf{1}\{w_d \neq 0\}$
- p -norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_p = (\sum_{d=1}^D |w_d|^p)^{\frac{1}{p}}$



Framework for supervised learning

$$y \in \mathbb{R} \quad \underline{h_{w,b}}(x) = w^T x + b$$

Squared loss

$$\sum_n |y_n - h_{w,b}(x_n)|$$

- Pick:

- ▶ Model/hypotheses ✓
- ▶ Loss function ✓
- ▶ Regularizer ✓
- ▶ Algorithm to solve optimization problem

These choices lead to different learning algorithms

Framework for machine learning

- Application
 - ▶ Labeled vs unlabeled data
 - ▶ Labeled: supervised learning. Type of label: categorical (classification), quantitative (regression)
 - ▶ Unlabeled: unsupervised learning.
- Model/hypotheses
- Optimization problem
- Algorithm to solve optimization problem

Summary

- Hypotheses/models: linear functions of features.
- Objective: choose a function (*i.e.*, parameters for the function) that minimize a cost function.
 - ▶ Cost function measures loss or error of the predictions made by a model/hypothesis on training set.
 - ▶ Cost function depends on the learning problem.
- Probabilistic interpretation
 - ▶ Minimizing cost function equivalent to maximizing likelihood.
 - ▶ Allows us to understand assumptions and to generalize our models.
- How do we minimize the cost function ?
 - ▶ Numerical methods: gradient and stochastic gradient descent.
 - ▶ Sometimes analytical solutions are available.
 - ▶ Computational considerations influence the choice.
- Can allow non-linear models/hypotheses by transforming features using non-linear functions.

Summary

- Regularization: preferring a simpler model by penalizing large weights or many non-zero weights.
- Helpful to reduce overfitting and to prevent numerical instability.
- Can be added to any linear model.
- Only a minor modification to the unregularized algorithms.
- New hyperparameter λ needs to be chosen using cross-validation.