

Machine Learning: The Optimization Perspective

Konstantin Tretyakov

http://kt.era.ee

AACIMP 2012 August 2012 STACC \$

Software Technology and Applications Competence Center









So far...

Machine learning is important and interesting

The general concept:

Fitting models to data



So far...

- Machine learning is important and interesting
- The general concept:







So far...







Optimization is important Optimization is possible





Optimization is important Optimization is possible*

* Basic techniques

- Constrained / Unconstrained
- Analytic / Iterative
- Continuous / Discrete



Special cases of optimization

Machine learning

••••



- Machine learning
- Algorithms and data structures
- General problem-solving
- Management and decision-making



Special cases of optimization

- Machine learning
- Algorithms and data structures
- General problem-solving
- Management and decision-making
- Evolution
- The Meaning of Life?





Optimization task

Given **a function**

$$f(\mathbf{x}):\mathbf{x}\to\mathbb{R}$$

find the argument x resulting in the optimal value.



Given a function

$$f(\mathbf{x}):\mathbf{x} \to \mathbb{R}$$

find the argument x resulting in the optimal value, subject to

$$\mathbf{x} \in \mathcal{C}$$



In principle, **x** can be anything:

Discrete

- Value (e.g. a name)
- Structure (e.g. a graph, plaintext)
- Finite / infinite

Continuous*

- Real-number, vector, matrix, …
- Complex-number, function, ...



In principle, **f** can be anything:

- Random oracle
- Structured
- Continuous
- Differentiable
- Convex







		Knowledge about f	
		Not much	A lot
Type of X	Discrete	Combinatorial search: Brute-force, Stepwise, MCMC, Population-based,	Algorithmic
	Continuous	Numeric methods: Gradient-based, Newton-like, MCMC, Population-based,	Analytic



	Knowledge	Knowledge about f	
Finding a weight-	Not much	A lot	
Type of Y	Combinatorial search: Brute-force, Stepwise, MCMC, Population-based,	Algorithmic	
Continuo	As A Section 2014 A S	Analytic	

D



		Knowledge about f	
Finding a weight- vector w, that minimizes the model error, in a fairly general case		Not much	A lot
		Combinatorial search: Brute-force, Stepwise, MCMC, Population-based,	Algorithmic
		Numeric methods: Gradient-based, Newton-like, MCMC, Population-based	Analytic



		Knowledge about f	
Finding a weight- vector w, that minimizes the model error, in a very general case		Not much	A lot
		Combinatorial search: Brute-force, Stepwise, MCMC, Population-based,	Algorithmic
		Numeric methods: Gradient-based, Newton-like, MCMC, Population-based,	Analytic

D



		Knowledge about f	
Finding a weight- vector w, that minimizes the model error, in many practical cases		Not much	A lot
		Combinatorial search: Brute-force, Stepwise, MCMC, Population-based,	Algorithmic
		Numeric methods: Gradient-based, Newton-like, MCMC, Population-based,	Analytic





Minima and maxima







Differentiability

D





Differentiability

D





Differentiability

Definition. We call a function $f : \mathbb{R}^m \to \mathbb{R}$ differentiable at point \mathbf{x}_0 if there exists $\mathbf{c}(\mathbf{x}_0) \in \mathbb{R}^m$ such that:

$$\Delta f(\mathbf{x}_0) = f(\mathbf{x}_0 + \Delta \mathbf{x}) - f(\mathbf{x}_0) = \mathbf{c}(\mathbf{x}_0)^T \Delta \mathbf{x} + o(\Delta \mathbf{x})$$

We call $\mathbf{c}(\mathbf{x}_0)$ the gradient or derivative^{*} of f (at point \mathbf{x}_0) and denote it by:

$$\frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}}$$
 or $f'(\mathbf{x}_0)$ or $\nabla f(\mathbf{x}_0)$



Let *f* be differentiable and let $\nabla f(\mathbf{x}_0) = \mathbf{c} \neq \mathbf{0}$. Take $\Delta \mathbf{x} = \dots$ Then:

$$f(\mathbf{x}_0 + \mathbf{\Delta}\mathbf{x}) \approx f(\mathbf{x}_0) + \mathbf{c}^T \dots < f(\mathbf{x}_0).$$

therefore \mathbf{x}_0 can't be a minimum of f.



Let *f* be differentiable and let $\nabla f(\mathbf{x}_0) = \mathbf{c} \neq \mathbf{0}$. Take $\Delta \mathbf{x} = -\mu \mathbf{c}$. Then:

$$f(\mathbf{x}_0 + \mathbf{\Delta}\mathbf{x}) \approx f(\mathbf{x}_0) + \mathbf{c}^T(-\mu\mathbf{c}) < f(\mathbf{x}_0).$$

therefore \mathbf{x}_0 can't be a minimum of f.



- This small observation gives us everything we need for now
 - A nice interpretation of the gradient
 - An extremality criterion
 - An **iterative algorithm** for function minimization



Interpretation of the gradient



D



Interpretation of the gradient



D



Theorem (Fermat). Let f be differentiable. Then

 \mathbf{x}_0 is an extremum $\Rightarrow \nabla f(\mathbf{x}_0) = \mathbf{0}$.

The converse does not hold in general.



Iterative algorithm

- I. Pick random point x_0
- 2. If $\nabla f(x_0) = 0$, then we've found an extremum.
- 3. Otherwise,



- I. Pick random point x_0
- 2. If $\nabla f(\mathbf{x}_0) = \mathbf{0}$, then we've found an extremum.
- 3. Otherwise, make a small step downhill:

$$\boldsymbol{x}_1 \leftarrow \boldsymbol{x}_0 - \mu_0 \nabla f(\boldsymbol{x}_0)$$



- I. Pick random point x_0
- 2. If $\nabla f(\mathbf{x}_0) = \mathbf{0}$, then we've found an extremum.
- 3. Otherwise, make a small step downhill: $\nabla f(x)$

$$\boldsymbol{x}_1 \leftarrow \boldsymbol{x}_0 - \mu_0 \nabla f(\boldsymbol{x}_0)$$

- 4. ... and then another step $x_2 \leftarrow x_1 \mu_1 \nabla f(x_1)$
- 5. ... and so on until



Gradient descent

- I. Pick random point x_0
- 2. If $\nabla f(\mathbf{x}_0) = \mathbf{0}$, then we've found an extremum.
- 3. Otherwise, make a small step downhill:

$$\boldsymbol{x}_1 \leftarrow \boldsymbol{x}_0 - \mu_0 \nabla f(\boldsymbol{x}_0)$$

- 4. ... and then another step $x_2 \leftarrow x_1 \mu_1 \nabla f(x_1)$
- 5. ... and so on until $\nabla f(\mathbf{x}_n) \approx \mathbf{0}$ or we're tired.

With a smart choice of μ_i we'll converge to a minimum



Gradient descent


Gradient descent

D



$\boldsymbol{x}_{i+1} \leftarrow \boldsymbol{x}_i - \mu_i \nabla f(\boldsymbol{x}_i)$

Gradient descent

Þ



$\Delta \boldsymbol{x}_i = -\mu_i \nabla f(\boldsymbol{x}_i)$

Gradient descent

Þ



$\Delta \boldsymbol{x}_i = -\mu \boldsymbol{c}$





 $\Delta \boldsymbol{x}_i = -\mu \, \nabla f(\boldsymbol{x}_i)$





 $\Delta \boldsymbol{x}_i = -\mu \, \nabla f(\boldsymbol{x}_i)$



D

Problem. Given a dataset $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbf{R}^m$, find \mathbf{w} , that minimizes

$$f(\mathbf{w}) = \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{w}\|^2$$

Propose an analytical as well as an iterative solution.



Þ



 x_1, x_2, \dots, x_{50}





l



























Stochastic gradient descent

Whenever

 $f(\boldsymbol{w}) = \sum g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$



$$f(\boldsymbol{w}) = \sum g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

e.g.

$$\operatorname{Error}(w) = \sum (\operatorname{model}_w(x_k) - y_k)^2$$



$$f(\boldsymbol{w}) = \sum g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

the gradient is also a sum:

$$\nabla f(\boldsymbol{w}) = \sum \nabla g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$



$$f(\boldsymbol{w}) = \sum g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

the gradient is also a sum:

$$\nabla f(\boldsymbol{w}) = \sum \nabla g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

The GD step is then also a sum

$$\Delta \boldsymbol{w}_{\boldsymbol{i}} = -\mu \sum \nabla g(\boldsymbol{w}_{\boldsymbol{i}}, \boldsymbol{x}_{\boldsymbol{k}})$$



$$f(\boldsymbol{w}) = \sum g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

the gradient is also a sum:

$$\nabla f(\boldsymbol{w}) = \sum \nabla g(\boldsymbol{w}, \boldsymbol{x}_{\boldsymbol{k}})$$

The GD step is then also a sum

$$\Delta \boldsymbol{w}_{\boldsymbol{i}} = -\mu \sum \nabla g(\boldsymbol{w}_{\boldsymbol{i}}, \boldsymbol{x}_{\boldsymbol{k}})$$



Batch update:

 $\Delta w_i = -\mu \sum \nabla g(w_i, x_k)$



Batch update:

$$\Delta \boldsymbol{w}_{\boldsymbol{i}} = -\mu \sum \nabla g(\boldsymbol{w}_{\boldsymbol{i}}, \boldsymbol{x}_{\boldsymbol{k}})$$

• On-line update:

$$\Delta \boldsymbol{w_i} = -\mu \, \nabla g(\boldsymbol{w_i}, \boldsymbol{x_{random}})$$









Summary

- An interpretation of the gradient
- An extremality criterion
- An iterative algorithm for function minimization
- A stochastic iterative algorithm for function minimization



D



• The symbol Δ is called _____ and denotes

• The symbol ${f V}$ is called _____ and denotes



Þ

Gradient descent:

$\Delta w_i =$



Þ

Gradient descent:

 $\Delta \boldsymbol{w}_i = -\mu \boldsymbol{\nabla}_{\boldsymbol{w}} \mathbf{f}(\boldsymbol{w})$



Þ

Gradient descent:

 $\Delta \boldsymbol{w}_i = -\mu \boldsymbol{c}$



Gradient descent:

$$\Delta \boldsymbol{w}_i = -\mu \boldsymbol{c}$$

Suppose the batch gradient descent step is $\Delta w_i = \sum_i (w^T x_i + |x_i|^2) e$ the corresponding stochastic gradient descent step is:

$$\Delta w_i =$$



Þ

Can bacteria learn?



Questions?





Linear Regression

Konstantin Tretyakov

http://kt.era.ee

AACIMP 2012 August 2012 STACC

Software Technology and Applications Competence Center







Supervised Learning

• Let X and Y be some sets.

• Let there be a training dataset: $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ $x_i \in X, y_i \in Y$



Supervised Learning

• Let X and Y be some sets.

• Let there be a training dataset: $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ $x_i \in X, y_i \in Y$

Supervised learning:

Find a function $f: X \rightarrow Y$, generalizing the dependency present in the data.



Classification



X = ℝ², Y = {blue, red} D = {((1.3, 0.8), red), ((2.5, 2.3), blue), ... } f(x₁, x₂) = if (x₁ + x₂) > 3 then blue else red



Regression





Linear Regression





Linear Regression

D

$X = \mathbb{R}^{m}, \qquad Y = \mathbb{R}$ $f(x_1, \dots, x_m) = w_0 + w_1 x_1 + \dots + w_m x_m$


$$X = \mathbb{R}^{m}, \qquad Y = \mathbb{R}$$
$$f(x_1, \dots, x_m) = w_0 + w_1 x_1 + \dots + w_m x_m$$

$$f(x_1, \dots, x_m) = w_0 + \langle w, x \rangle$$

Inner product



D

$$X = \mathbb{R}^{m}, \qquad Y = \mathbb{R}$$
$$f(x_1, \dots, x_m) = w_0 + w_1 x_1 + \dots + w_m x_m$$

$$f(x_1, \dots, x_m) = w_0 + \langle \boldsymbol{w}, \boldsymbol{x} \rangle$$
$$f(x_1, \dots, x_m) = w_0 + (w_1, \dots, w_m) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$



Þ

$$X = \mathbb{R}^{m}, \qquad Y = \mathbb{R}$$
$$f(x_1, \dots, x_m) = w_0 + w_1 x_1 + \dots + w_m x_m$$

$$f(x_1, \dots, x_m) = w_0 + \langle \boldsymbol{w}, \boldsymbol{x} \rangle$$
$$f(x_1, \dots, x_m) = w_0 + (w_1, \dots, w_m) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$
$$f(x_1, \dots, x_m) = w_0 + \boldsymbol{w}^T \boldsymbol{x}$$



 $f(\boldsymbol{x}) = w_0 + \boldsymbol{w}^T \boldsymbol{x}$



D

$$f(\mathbf{x}) = w_0 + \mathbf{w}^T \mathbf{x}$$

Bias term



$$f(\boldsymbol{x}) = w_0 + \boldsymbol{w}^T \boldsymbol{x}$$

$$f(x_1, \dots, x_m) = (w_0, w_1, \dots, w_m) \begin{pmatrix} 1\\ x_1\\ \vdots\\ x_m \end{pmatrix}$$



D

$$f(\boldsymbol{x}) = w_0 + \boldsymbol{w}^T \boldsymbol{x}$$

$$f(x_1, \dots, x_m) = (w_0, w_1, \dots, w_m) \begin{pmatrix} 1\\ x_1\\ \vdots\\ x_m \end{pmatrix}$$

$$f(\boldsymbol{x}) = \widetilde{\boldsymbol{w}}^T \widetilde{\boldsymbol{x}}$$



Þ

$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x}$











Þ





Þ





Þ



August, 2012







Þ

 $E_D(\boldsymbol{w}) = \sum_i (\boldsymbol{w}^T \boldsymbol{x}_i - \boldsymbol{y}_i)^2$



$$E_D(w) = \sum_i (x_i^T w - y_i)^2$$
$$\begin{pmatrix} \vdots \\ w \\ \vdots \end{pmatrix}$$
$$(\dots x_1^T \dots) \quad x_1^T w \qquad y_1$$
$$(\dots x_2^T \dots) \quad x_2^T w \qquad y_2$$
$$\vdots$$
$$(\dots x_n^T \dots) \quad x_n^T w \qquad y_n$$







Þ

 $E_D(\boldsymbol{w}) = \sum_i (\boldsymbol{x}_i^T \boldsymbol{w} - \boldsymbol{y}_i)^2$

Xw - v



Þ

 $E_D(\boldsymbol{w}) = \sum_{i} (\boldsymbol{x}_i^T \boldsymbol{w} - \boldsymbol{y}_i)^2$

 $||Xw - y||^2$



Þ

 $E_D(\boldsymbol{w}) = \sum_i (\boldsymbol{x}_i^T \boldsymbol{w} - \boldsymbol{y}_i)^2$

 $||Xw - y||^2$

 $Xw \approx y$



D

$$E_D(\boldsymbol{w}) = \sum_i (\boldsymbol{x}_i^T \boldsymbol{w} - \boldsymbol{y}_i)^2$$

 $||Xw - y||^2$

 $Xw \approx y$

AACIMP Summer School. August, 2012

 $w \approx X^{-1}y?$



D

$$E_D(\boldsymbol{w}) = \sum_i (\boldsymbol{x}_i^T \boldsymbol{w} - \boldsymbol{y}_i)^2$$

 $||Xw - y||^2$

 $w = X^+ y !$ $Xw \approx y$



 $\operatorname{argmin}_{\mathbf{w}} \| X \mathbf{w} - \mathbf{y} \|^2$



Þ

 $\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \| \mathbf{X}\mathbf{w} - \mathbf{y} \|^2$



$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \| \mathbf{X}\mathbf{w} - \mathbf{y} \|^{2}$$
$$\frac{\| \mathbf{a} \|^{2} = \mathbf{a}^{T} \mathbf{a}}{\frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^{T} (\mathbf{X}\mathbf{w} - \mathbf{y})}$$

.



$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \| \mathbf{X}\mathbf{w} - \mathbf{y} \|^2$$

$$\frac{1}{2}(Xw - y)^T(Xw - y)$$
$$(a + b)^T = (a^T + b^T)$$
$$\frac{1}{2}(w^TX^T - y^T)(Xw - y)$$



1

Þ

 $\frac{1}{2}(\boldsymbol{w}^T\boldsymbol{X}^T - \boldsymbol{y}^T)(\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})$

a(b+c) = ab + ac

$$\frac{1}{2} \left(\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} - \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{y}^T \boldsymbol{y} \right)$$



Þ

 $\frac{1}{2}(\boldsymbol{w}^T\boldsymbol{X}^T - \boldsymbol{y}^T)(\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})$

$$\frac{1}{2} \left(\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} - \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{y}^T \boldsymbol{y} \right)$$
$$\boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} = \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} = \text{scalar}$$
$$\frac{1}{2} \left(\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - 2 \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{y}^T \boldsymbol{y} \right)$$



D

 $E(\boldsymbol{w}) = \frac{1}{2} (\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - 2\boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{y}^T \boldsymbol{y})$



$$E(w) = \frac{1}{2} (w^T X^T X w - 2y^T X w + y^T y)$$
$$\nabla(f + g) = \nabla f + \nabla g$$
$$\nabla(w^T A w) = 2Aw$$
$$\nabla(a^T w) = a$$

$$\nabla E(\boldsymbol{w}) = \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - \boldsymbol{X}^T \boldsymbol{y}$$



D

 $E(\boldsymbol{w}) = \frac{1}{2} (\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - 2\boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{y}^T \boldsymbol{y})$

$$\nabla E(w) = X^T X w - X^T y$$

$$\mathbf{0} = X^T X w - X^T y$$

$$X^T X w = X^T y$$



Þ

 $E(\boldsymbol{w}) = \frac{1}{2} (\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} - 2\boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{y}^T \boldsymbol{y})$

$$\nabla E(w) = X^{T}Xw - X^{T}y$$
$$\mathbf{0} = X^{T}Xw - X^{T}y$$
$$X^{T}Xw = X^{T}y$$

$$\boldsymbol{w} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$



D

 $\boldsymbol{w} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$



 $\boldsymbol{w} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$

X = matrix(X) y = matrix(y) w = (X.T * X).I * X.T * y



$$\boldsymbol{w} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

Moore-Penrose pseudoinverse

- X = matrix(X)
- y = matrix(y)
- w = (X.T * X).I * X.T * y



D

$$w = (X^{T}X)^{-1}X^{T}y$$

$$w = X^{+}y$$
Moore-Penrose
pseudoinverse
$$y = matrix(X)$$

$$w = (X.T * X).I * X.T * y$$

$$w = pinv(X) * y$$


from sklearn.linear_model import
 LinearRegression

model = LinearRegression()
model.fit(X, y)

w=(model.intercept_,model.coef_)
model.predict(X_new)



Stochastic Gradient Regression

 $\Delta \boldsymbol{w} = -\mu(\boldsymbol{w}^T\boldsymbol{x}_i - \boldsymbol{y}_i)\boldsymbol{x}_i$



Stochastic Gradient Regression

$\Delta w = -\mu e_i x_i$



Stochastic Gradient Regression

$\Delta w = -\mu e_i x_i$

from sklearn.linear model import SGDRegressor

model = SGDRegressor(alpha=0, n_iter=30) model.fit(X, y)



Polynomial Regression

Say we'd like to fit a model:

 $f(x_1, x_2)$ $= w_0 + w_1 x_1 + w_2 x_2^2 + w_3 x_1 x_2$



Polynomial Regression

Say we'd like to fit a model:

$$f(x_1, x_2) = w_0 + w_1 x_1 + w_2 x_2^2 + w_3 x_1 x_2$$

Simply transform the features and proceed as normal:

$$(x_1, x_2) \rightarrow (x_1, x_2^2, x_1 x_2)$$



- # n x 1 matrix
- x = matrix(...)

- # Add bias & square features
 X = hstack([x**0, x**1, x**2])
- # Solve for w
 w = pinv(X) * y



Overfitting

D





D





D





$$E(\boldsymbol{w}) = \frac{1}{2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 + \lambda \|\boldsymbol{w}\|_1$$
$$\ell_2 \text{-loss} \qquad \ell_1 \text{-penalty}$$



Regularization

$$E(\boldsymbol{w}) = \frac{1}{2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 + \lambda \|\boldsymbol{w}\|_0$$
$$\ell_2 \text{-loss} \qquad \ell_0 \text{-penalty}$$



Regularization

$$E(\boldsymbol{w}) = \frac{1}{2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 + \lambda \|\boldsymbol{w}\|_0$$
$$\ell_2 \text{-loss} \qquad \ell_0 \text{-penalty}$$

>>> SGDRegressor?

Parameters

loss : str, 'squared_loss' or 'huber' ...
...
penalty : str, '12' or '11' or 'elasticnet'



D





Ridge regression

$$\operatorname{argmin}_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 + \lambda \|\boldsymbol{w}\|^2$$

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





Ridge regression

$$\operatorname{argmin}_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 + \lambda \|\boldsymbol{w}_*\|^2$$

$$\boldsymbol{w} = (\boldsymbol{X}^T \boldsymbol{X} + \lambda \boldsymbol{I}_*)^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

The bias term w_0 is usually **not penalized**.







Derive an SGD algorithm for Ridge Regression.







OLS linear regression searches for a model that has the best

uiz

• Analytic solution for OLS regression: w =_____

Stochastic gradient solution for OLS regression:

$$\Delta w =$$



Large number of model parameters and/or small data may lead to _____.

We address overfitting by _

uiz

 "Ridge regression" means ___-loss and ___penalty.

Analytic solution for Ridge regression:





As we increase regularization strength (i.e. increase λ), the training error _____.

... and the test error _____



Questions?

