

Machine Learning: The Optimization Perspective

Konstantin Tretyakov

http://kt.era.ee

IFI Summer School 2014 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...





The Land of Machine Learning





The Land of Machine Learning





The Land of Machine Learning



What do you need to know about optimization?

Þ



What do you need to know about optimization?



Optimization is important Optimization is possible

What do you need to know about optimization?



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



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





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.





 x_1, x_2, \dots, x_{50}















June 2014















 $\Delta \mathbf{w} = \mu \cdot 2(\mathbf{x}_i - \mathbf{w})$



Konstantin Tretyakov

http://kt.era.ee

IFISS Summmer School 2014 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









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



Þ

$$X = \mathbb{R}^{m}, \qquad Y = \mathbb{R}$$
$$f(x_{1}, \dots, x_{m}) = w_{0} + w_{1}x_{1} + \dots + w_{m}x_{m}$$
$$f(x_{1}, \dots, x_{m}) = w_{0} + \langle w, x \rangle$$

Inner product



$$X = \mathbb{R}^{m}, \quad Y = \mathbb{R}$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + w_{1}x_{1} + \dots + w_{m}x_{m}$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + \langle w, x \rangle$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + (w_{1}, \dots, w_{m}) \begin{pmatrix} x_{1} \\ x_{2} \\ \vdots \\ x_{m} \end{pmatrix}$$



Þ

$$X = \mathbb{R}^{m}, \quad Y = \mathbb{R}$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + w_{1}x_{1} + \dots + w_{m}x_{m}$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + \langle w, x \rangle$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + (w_{1}, \dots, w_{m}) \begin{pmatrix} x_{1} \\ x_{2} \\ \vdots \\ x_{m} \end{pmatrix}$$

$$f(x_{1}, \dots, x_{m}) = w_{0} + w^{T}x$$



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



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

Bias term



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

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



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

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

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


Linear Regression

Þ

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



Linear Regression

























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



D









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

Xw - v



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

 $||Xw - y||^2$



Þ

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

 $||Xw - y||^2$

 $Xw \approx y$



Þ

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

 $||Xw - y||^2$

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



Þ

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

 $||Xw - y||^2$

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



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



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



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





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

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

$$(a + b)^T = (a^T + b^T)$$



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

a(b+c) = ab + ac

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



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

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



Þ

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



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

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



Þ

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

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

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

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



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

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

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

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

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



Þ

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



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

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



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

$$w = X^{+}y$$
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

D

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



Stochastic Gradient Regression

D

 $\Delta w = -\mu e_i x_i$



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:

D

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



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:

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



Single-variable Polynomial OLS

- # 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













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



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'







Ridge regression

D

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

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





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.







Quiz

 OLS linear regression searches for a _ model that has the best _____.

• 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:

$$w = _$$





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

... and the test error _____



Questions?

