Guided Tour of Machine Learning in Finance

Linear Regression

Igor Halperin

NYU Tandon School of Engineering, 2017

Task: predict a scalar value $y \in \mathbb{R}$ from a vector of predictors ("**features**") $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$ ($D \ge 1$ is the dimension of the feature space, or the number of predictors).

Task: predict a scalar value $y \in \mathbb{R}$ from a vector of predictors ("**features**") $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$ ($D \ge 1$ is the dimension of the feature space, or the number of predictors).

Given: a dataset
$$(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$$

Task: predict a scalar value $y \in \mathbb{R}$ from a vector of predictors ("**features**") $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$ ($D \ge 1$ is the dimension of the feature space, or the number of predictors).

Given: a dataset
$$(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$$

Architecture: Linear

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{W}$$

N predicted values

X
is a N x D
design matrix

 $\mathbf{W} \in \mathbb{R}^D$ is a vector of parameters

Task: predict a scalar value $y \in \mathbb{R}$ from a vector of predictors ("**features**") $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$ ($D \ge 1$ is the dimension of the feature space, or the number of predictors).

Given: a dataset
$$(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$$

Architecture: Linear

Example: y is a vector of N daily returns of AMZN, and **X** is a $N \times D$ design matrix made of D-1 daily market index returns (S&P 500, NASDAQ, VIX, etc.)

Performance measure P: mean square error (MSE) on the test set:

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2$$

Performance measure P: mean square error (MSE) on the test set:

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} ||\hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test}||_2^2$$

Performance measure P: mean square error (MSE) on the test set:

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_i^{test} - y_i^{test} \right)^2 = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_2^2 = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_2^2$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{\mathbf{y}}_{i}^{test} - \mathbf{y}_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$?

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_{i}^{test} - y_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$?

Minimize MSE on the training set!

Why? Because both MSE_{train} and MSE_{test} estimate the same generalization (expected) error $\mathbb{E}\left[\left(\hat{\mathbf{y}}-\mathbf{y}\right)^2\right]$ from the empirical distribution $\sim p_{data}$.

Note: In ML, we often minimize one function, while actually caring about minimization of another function. This makes ML different from optimization that typically focuses only on one function.

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_{i}^{test} - y_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

$$\nabla_{W} \mathbf{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left| \left| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right| \right|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left| \left| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right| \right|_{2}^{2} = 0$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{\mathbf{y}}_{i}^{test} - \mathbf{y}_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

$$\nabla_{W} \text{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_{i}^{test} - y_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

$$\nabla_{W} \mathbf{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

$$\mathbf{X}^{train} \rightarrow \mathbf{X}, \mathbf{y}^{train} \rightarrow \mathbf{y}$$
Simplify notation:
$$\mathbf{X}^{train} \rightarrow \mathbf{X}, \mathbf{y}^{train} \rightarrow \mathbf{y}$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{\mathbf{y}}_{i}^{test} - \mathbf{y}_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

$$\nabla_{W} \text{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right)$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right)$$
Simplify notation:
$$\mathbf{X}^{train} \rightarrow \mathbf{X}, \mathbf{y}^{train} \rightarrow \mathbf{y}$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_{i}^{test} - y_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

$$\nabla_{W} \text{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_{i}^{test} - y_{i}^{test} \right)^{2} = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_{2}^{2} = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_{2}^{2}$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

For optimal weights \mathbf{W} , the gradient of \mathbf{MSE}_{train} should be 0:

$$\nabla_{W} \text{MSE}_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$
Simplify notation:
$$\mathbf{X}^{train} \rightarrow \mathbf{X}, \mathbf{y}^{train} \rightarrow \mathbf{y}$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$

$$\Rightarrow \qquad \mathbf{W} = \left(\mathbf{X}^T \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{y}$$

Normal equation

Performance measure P: mean square error (MSE) on the test set:

$$\mathbf{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} \left(\hat{y}_i^{test} - y_i^{test} \right)^2 = \frac{1}{N_{test}} \left| \left| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right| \right|_2^2 = \frac{1}{N_{test}} \left| \left| \mathbf{X}^{test} \cdot \mathbf{W} - \mathbf{Y}^{test} \right| \right|_2^2$$

How to optimize parameters $\mathbf{W} \in \mathbb{R}^D$? Minimize MSE on the **training set**! Why? Because both \mathbf{MSE}_{test} and \mathbf{MSE}_{train} estimate the generalization (expected) error.

For optimal weights \mathbf{W} , the gradient of \mathbf{MSE}_{train} should be 0:

$$\nabla_{W} MSE_{train} = \nabla_{W} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_{2}^{2} = \frac{1}{N_{train}} \nabla_{W} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_{2}^{2} = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right)^{T} \left(\mathbf{X}^{train} \mathbf{W} - \mathbf{y}^{train} \right) = 0$$

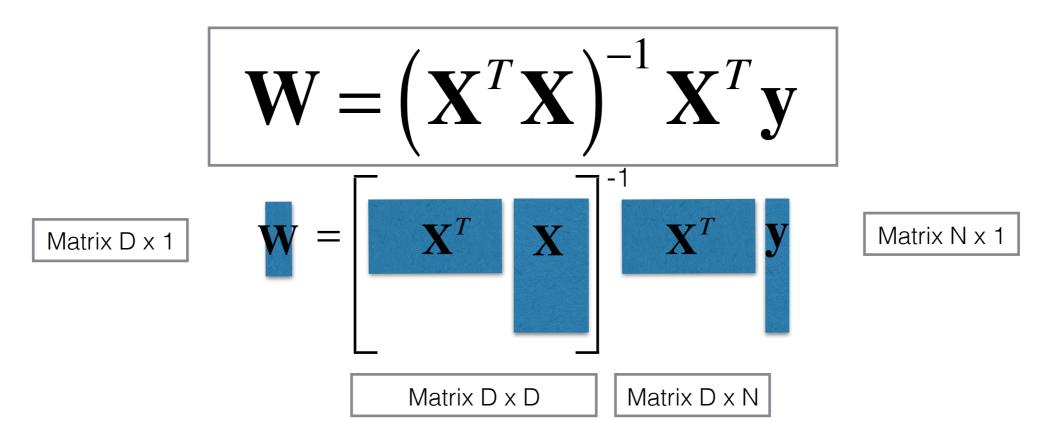
$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$

$$\Rightarrow \nabla_{W} \left(\mathbf{W}^{T} \mathbf{X}^{T} \mathbf{X} \mathbf{W} - 2 \mathbf{W}^{T} \mathbf{X}^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y} \right) = 2 \left(\mathbf{X}^{T} \mathbf{X} \right) \mathbf{W} - 2 \left(\mathbf{X}^{T} \mathbf{y} \right) = 0$$

$$\mathbf{X} = \mathbf{X}^{train}, \mathbf{y} = \mathbf{y}^{train} \Rightarrow \mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Normal equation

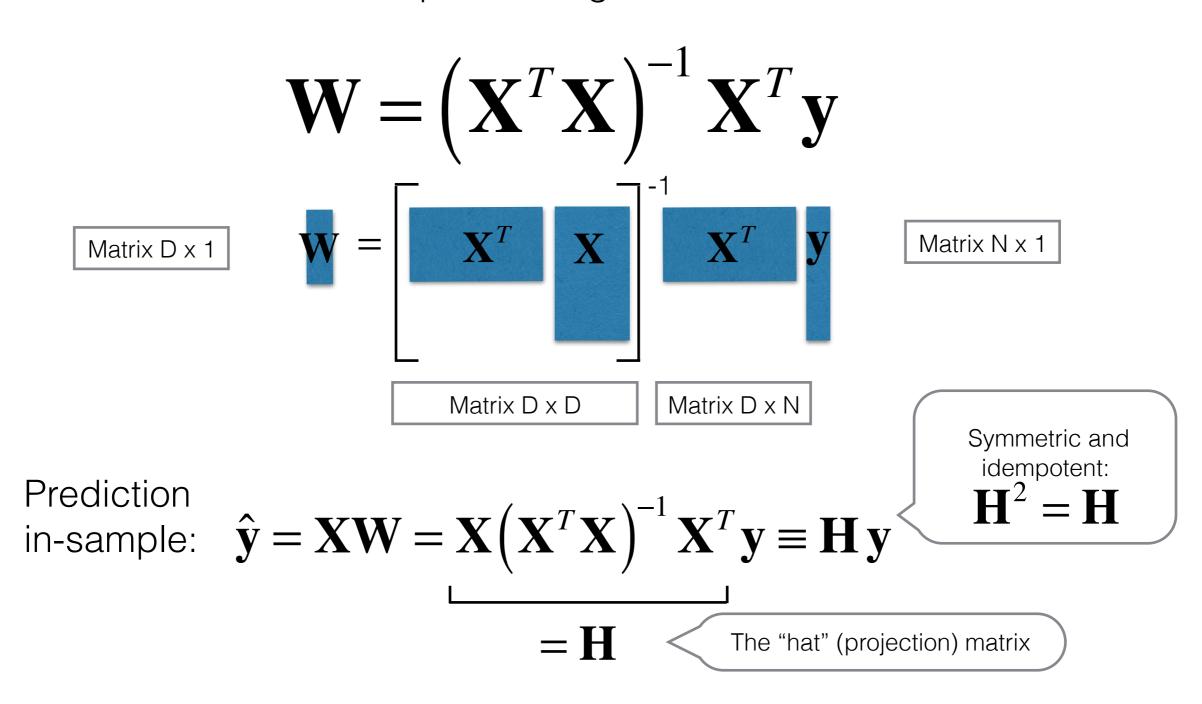
Visualize the result for optimal weights W in the vector form



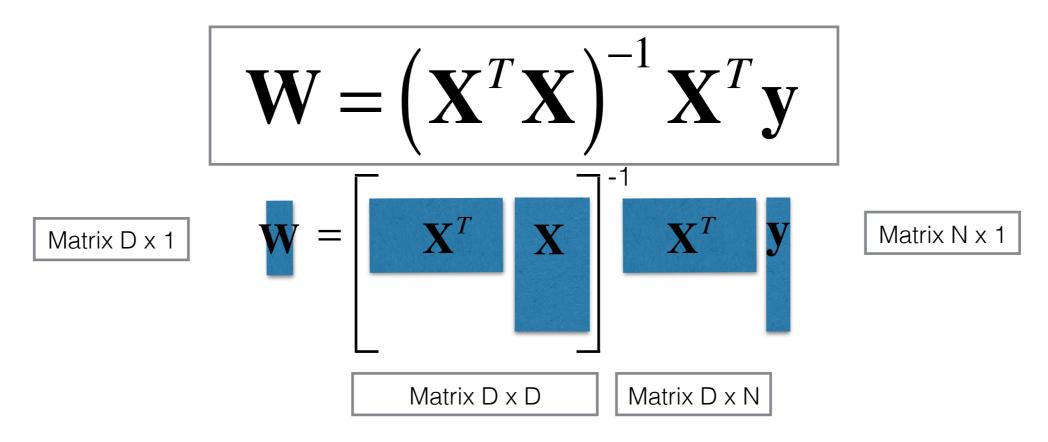
Prediction in-sample:
$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{W} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \equiv \mathbf{H} \mathbf{y}$$

$$= \mathbf{H} \qquad \text{The "hat" (projection) matrix}$$

Visualize the result for optimal weights W in the vector form



Visualize the result for optimal weights W in the vector form



Prediction out-of-sample: $\hat{\mathbf{y}} = \mathbf{X}^{test}\mathbf{W}$

How to overfit with Linear Regression: add more features to the design matrix without a proper control!