

TAM 598

Lecture 13 :

Supervised Learning - Linear Regression  
via Least Squares

---

Announcements:

- HW 3 covers lectures 8-12 ; due on Mar 12
- One recorded class over spring break
- No class Monday Mar 24<sup>th</sup>

Supervised learning - given pairs of  $x$  and  $y$  observations,  
find the map between  $x$  and  $y$ .

**Classification** - When  $y$  is discrete  $\rightarrow$  aka **logistic regression**  
**Regression** - When  $y$  is continuous

## I. Linear Regression via Least Squares - Supervised Learning

given  $n$  d-dimensional inputs  $\underline{x}_{1:n} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n\}$   $\nwarrow$  features

and  $n$  outputs :  $y_{1:n} = \{y_1, \dots, y_n\}$   $\nwarrow$  targets

"linear" means that the map from  $\underline{x}$  to  $y$  is linear

Two common ways to train a supervised learning model

(1) **optimization** - find parameters by maximizing the likelihood  
of the data (or minimize the neg. log likelihood, the loss)

(2) **probabilistic inference** : (a) Bayes to get posterior given data,

(b) use posterior to make predictions

## II. Linear Regression with a single variable - Example

a synthetic data set  $x, y :$

$$y_i = -0.5 + 2x_i + 0.1 \epsilon_i$$

$$\begin{array}{c} \downarrow \\ \text{input} \\ x_i \sim U([0, 1]) \end{array}$$

$\hookrightarrow$  gaussian noise  
 $\epsilon_i \sim N(0, 1)$

use least squares to fit to the linear model

$$y = w_0 + w_1 x \quad \tilde{w} = (w_0, w_1)^\top \text{ parameters}$$

by minimizing the square loss

$$L(\tilde{w}) = \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2 = \|\tilde{y} - \tilde{x} \tilde{w}\|^2$$

$$L(\underline{w}) = \sum_{i=1}^n (y_i - w_0 - \underline{w}_i x_i)^2 = \|\underline{y} - \underline{\underline{x}} \underline{w}\|^2$$

$\underline{y}$  = vector of observations =  $(y_1 \dots y_n)$

$\underline{w}$  = weight vector =  $(w_0, \underline{w}_i)^T$

$\underline{\underline{x}}$  = design matrix =  $\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$

LHS

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix} = \begin{bmatrix} 4 & x_1 + x_2 + x_3 + x_4 \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$n \times 2$

$\boxed{\begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}}$

RHS

$$= \begin{bmatrix} \underline{E}_{\underline{y}\underline{y}} \\ \underline{E}_{\underline{y}\underline{x}} \end{bmatrix}$$

$$\underline{\underline{x}}^T \underline{\underline{x}} \underline{w} = \underline{\underline{x}}^T \underline{y}$$

$$L = \|\underline{r}\|^2 = \underline{r}^T \underline{r}$$

to minimize:

$$\nabla_{\underline{w}} L = 2 \|\underline{r}\| \frac{d \|\underline{r}\|}{d \underline{w}} = 0$$

Note:  $\frac{d \|\underline{r}\|}{d r} = \text{sign}(\underline{r})$

$$\begin{aligned} \text{so } \frac{d L}{d \underline{w}} &= 2 \|\underline{r}\| \frac{d \|\underline{r}\|}{d \underline{w}} \\ &= 2 \|\underline{r}\| \frac{d \|\underline{r}\|}{d r} \frac{dr}{d \underline{w}} \\ &= 2 \|\underline{r}\| \text{sign}(\underline{r}) (-\underline{x}) \\ &= 2 (\underline{y} - \underline{\underline{x}} \underline{w}) (-\underline{x}) \end{aligned}$$

$$\Rightarrow 0 = \underline{y} - \underline{\underline{x}} \underline{w}$$

$$\underline{\underline{x}} \underline{w} = \underline{y}$$

$\boxed{\underline{\underline{x}}^T \underline{\underline{x}} \underline{w} = \underline{\underline{x}}^T \underline{y}}$  solution

Example 2: how things can go wrong

try instead:  $y_i = -0.5 + 2x_i + 2x_i^2 + \epsilon_i$

and we still try to fit a linear model.

This is **underfitting**

- model is too simple to capture underlying patterns.
- often happens when model has high bias, performs poorly on training & unseen test data

### III. Polynomial Regression

$$y_i = -0.5 + 2x_i + 2x_i^2 + \epsilon_i$$

Now we try  $y = w_0 + w_1 x + w_2 x^2$  and minimize the square loss

$$L(w_0, w_1, w_2) = \sum_{i=1}^n (y_i - w_0 - w_1 x_i - w_2 x_i^2)^2$$

$$L(\underline{w}) = \| \underline{y} - \underline{\underline{x}} \underline{w} \|^2$$

$$\text{where } \underline{w} = (w_0, w_1, w_2)^T$$

$$\underline{\underline{x}} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}$$

$$\underline{y} = (y_1, \dots, y_n)^T \quad (n \times 1)$$

$$\left. \begin{array}{l} (3 \times 1) \\ (n \times 3) \end{array} \right\} \text{another linear system}$$

$$\underline{\underline{x}}^T \underline{\underline{x}} \underline{w} = \underline{\underline{x}}^T \underline{y}$$

solve for  $\underline{w}$   
linear in  $\underline{w}$

Example 2 : what if we used a higher degree polynomial?

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^e \\ 1 & x_2 & x_2^2 & \dots & x_2^e \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^e \end{bmatrix} \quad (n \times p)$$

$$\hat{w} = (w_0, w_1, \dots, w_p)^T \quad (p \times 1)$$

if  $p$  becomes too large : **overfitting**

model is too complex relative to amount of training data

model learns not only the true patterns, but also  
noise & fluctuations.

performs very well on training data, but poorly  
on new, unseen data

## IV The Generalized Linear Model

$$y(\underline{x}; \underline{w}) = \sum_{j=1}^m w_j \phi_j(\underline{x}) = \underline{w}^\top \underline{\phi}(\underline{x})$$

weights  $\underline{w}$        $\underline{\phi} = (\phi_1, \dots, \phi_m)^\top$   
arbitrary basis functions

model is linear in  $\underline{w}$ , but basis functions  $\phi(\underline{x})$  can be non-linear

e.g.) our polynomial model :  $\phi_1(x) = 1$      $\phi_2(x) = x$      $\phi_3(x) = x^2$

e.g.) multivariate linear regression where  $\underline{x} = (x_1, \dots, x_d)$

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

is also a gen. linear model where  $\phi_1(\underline{x}) = 1$

$$\phi_2(\underline{x}) = x_1$$

$$\phi_3(\underline{x}) = x_2$$

⋮

eg) multi-dimensional polynomials

$$\phi_j(\tilde{x}) = \sum_{\alpha \in A_j} \beta_\alpha \tilde{x}^\alpha$$

eg) radial basis functions

$$\phi_j(\tilde{x}) = \exp \left\{ - \frac{\|\tilde{x} - \tilde{x}_j\|^2}{2\ell^2} \right\}$$

eg) Fourier series

$$\phi_{2j}(x) = \cos \left( \frac{2\pi j}{L} x \right)$$

$$\phi_{2j+1}(x) = \sin \left( \frac{2\pi j}{L} x \right)$$

How to fit the generalized model using least squares:

quadratic loss function

$$\begin{aligned} L(\underline{\omega}) &= \sum_{i=1}^N [y_i(x_i; \underline{\omega}) - y_i]^2 \\ &= \| \underline{\Phi} \underline{\omega} - \underline{y} \|^2 \\ &= (\underline{\Phi} \underline{\omega} - \underline{y})^T (\underline{\Phi} \underline{\omega} - \underline{y}) \end{aligned}$$

$\underline{\Phi}$  = design matrix       $\underline{\Phi} \in \mathbb{R}^{N \times M}$ ,     $\underline{\Phi}_{ij} = \phi_j(x_i)$

N = # observations ;    M = # basis functions

To minimize the loss :

① Take  $\nabla_{\underline{w}} \mathcal{L}(\underline{w})$ , set = 0, solve for  $\underline{w}$

solution :  $(\begin{smallmatrix} \underline{\Phi}^T & \underline{\Phi} \\ \underline{\underline{y}} & \underline{\underline{y}} \end{smallmatrix}) \underline{w} = (\begin{smallmatrix} \underline{\Phi}^T \\ \underline{\underline{y}} \end{smallmatrix}) \underline{y}$  (as before)

{ solve with numpy.linalg.lstsq  
provide  $\underline{\Phi}$ ,  $\underline{y}$  and it returns  $\underline{w}$

## IV. Measures of Predictive Accuracy

### (1) Training and test datasets

- not a good idea to test how good your model is using the training dataset
- need to test model on unseen data - a test data set
- often, take your data and split (eg 70% train, 30% test)
- can plot predictions vs. observations

### (2) Mean-squared error - scalar measure of goodness of fit

$$\text{MSE} = \frac{1}{n_t} \sum_{i=1}^{n_t} [y_{t,i} - w^\top \phi(x_{t,i})]^2$$

hard to understand the absolute meaning of MSE, so we can instead use relative MSE (RMSE)

$$\text{RMSE} = \frac{\text{MSE of your model}}{\text{MSE of simplest possible model}}$$

Simplest:  $y_{\text{simplest}} = \hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$

if  $\text{RMSE} < 1 \rightarrow$  you are doing better than the simplest possible model

$\text{RMSE} > 1 \rightarrow$  the simplest possible model outperforms your

### (3) Coefficient of Determination $R^2$

$R^2 = 1 - \text{RMSE}$  : tells us what % of the variance of the test data is explained by your model