

Regression

• Classification: given pt. x , predict class (discrete)

• Regression: given pt x , predict a numerical value

- choose:

① form of regression function $h(x; w)$ w/ parameters w
↳ like decision function from classification

② cost function to optimize

↳ usually based on loss function, e.g. empirical risk = $E[\text{loss}]$

Some regression fns:

① linear regression $h(x; w, \alpha) = w \cdot x + \alpha$

② polynomial regression (plug in polynomial features into linear regression)

③ logistic regression $h(x; w, \alpha) = S(w \cdot x + \alpha)$, $S(\cdot)$ is Sigmoid, $S(r) = \frac{1}{1 + e^{-r}}$

↳ LDA produces logistic posterior probabilities

↳ interpolate probabilities, not numbers

Loss Functions: z = prediction $h(x)$, y is label

① Squared error: $(z - y)^2$

② Absolute error: $|z - y|$ computationally harder to do

③ logistic/cross-entropy: $-y \ln z - (1 - y) \ln(1 - z)$ z & y must be in same range

Some Cost fns

① $J(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i)$ mean loss

② $J(h) = \max_{i=1}^n L(h(x_i), y_i)$ maximum loss

③ $J(h) = \sum_{i=1}^n w_i \cdot L(h(x_i), y_i)$ weighted sum

④ L2 regularization: $J(h) + \lambda \|w\|_2^2$

⑤ L1 regression: $J(h) + \lambda \|w\|_1$

} don't trust large weights

Famous Regression methods:

Least-Squares linear regression: ① + ④ + ⑤

Weighted LS linear: ① + ④ + ⑤

Ridge Regression: ① + ④ + ④

Lasso: ① + ④ + ⑤

Logistic Regression: ③ + ③ + ③

Least absolute deviations: ① + ② + ③

Chebyshev criterion: ① + ② + ③

} quadratic cost fn. (convex)

quadratic program (like SVMs)

convex cost; minimize w/ gradient descent

} linear programs

Least-Squares Linear Regression Gauss, 1801

Linear regression fn ① + Squared loss (A) + cost fn (C)

problem: find w, α that minimizes $\sum_{i=1}^n (\underbrace{x_i \cdot w + \alpha}_{\text{prediction}} - y_i)^2$

design matrix:

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ \vdots & \vdots & & \vdots & \vdots \\ x_{n1} & \dots & \dots & x_{nd} & 1 \end{bmatrix} \leftarrow \text{point } x_i^T \quad \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

↑ feature column x_{*d}

Usually, $n > d$ (tall, not wide)

Recall fictitious dimension trick: $h(x) = xw + \alpha$ as

$$\begin{bmatrix} x_1 & \dots & x_d & 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ d \end{bmatrix}$$

Now $x \in \mathbb{R}^{n \times (d+1)}$ w is a $d+1$ vector

find w that minimizes $\|xw - y\|^2 = \text{RSS}(w)$, residual sum of squares

Optimize w/ calculus:

$$\min \text{RSS}(w) = w^T x^T x w - 2y^T x w + y^T y$$

$$\nabla \text{RSS} = 2x^T x w - 2x^T y = 0$$

$$\Rightarrow x^T x w = x^T y \quad (\text{normal equations})$$

$$\Rightarrow w^* = (x^T x)^{-1} x^T y$$

at least
PSD

$$\Rightarrow w^* = \underbrace{(x^T x)^{-1}}_{\text{pseudo inverse of } x} x^T y \quad x^T \in \mathbb{R}^{(d+1) \times n}$$

If $x^T x$ is singular, this problem is underconstrained

↳ all points lie on a hyperplane, don't use all $d+1$ dimensions

We can use a linear solver to get w^*

$$\text{Observe: } x^T x = (x^T x)^{-1} x^T x = I_{d+1}$$

Observe: predicted value of y_i is $\hat{y}_i = w \cdot x_i$

What are y_i ?

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = XW = XX^T y = Hy \quad H \in \mathbb{R}^{n \times n}$$

$H = XX^T$ is the **hat matrix**

If $n > d+1$, H is Singular \implies no perfect fit w/ hyperplane

Advantages of LS LR:

- easy to compute, just solve linear system
- usually unique & stable solution

Disadvantages

- Very sensitive to outliers since errors are squared
- fails if XX^T is singular

Logistic Regression

David Cox, 1958

$$h = S(x; w+d)$$

$$-y \ln z - (1-y) \ln(1-z)$$

$$\frac{1}{n} \sum L(h(x_i), y_i)$$

logistic regression fn ③ + logistic loss fn ② + cost fn ①

- fits probabilities in range $(0, 1)$
- usually used for classification
 - \hookrightarrow input y_i 's might be probabilities, in practice $y_i \in \{0, 1\}$

- QDA & LDA are generative models used to calculate posterior probabilities

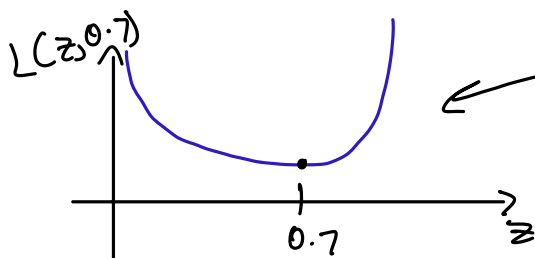
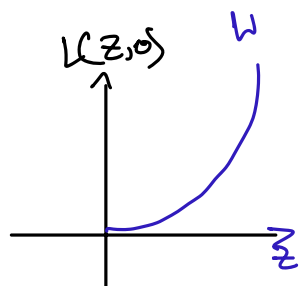
- logistic regression: **discriminative model**

\hookrightarrow build model of what we want to predict

\hookrightarrow no need to fit each class

w/ X and W (including fictitious dimension; d is w 's last component)

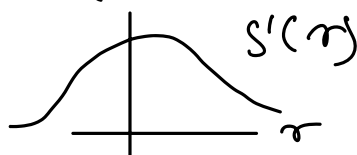
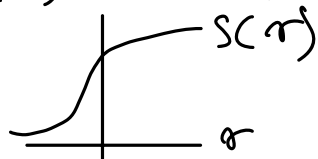
$$\text{find } w^* = \min J = \min \sum_{i=1}^n L(S(x_i \cdot w), y_i) = - \sum_{i=1}^n (y_i \ln s(x_i \cdot w) + (1-y_i) \ln(1-s(x_i \cdot w)))$$



$\leftarrow y_i \notin \{0, 1\}$

$J(w)$ is **convex!** \rightarrow use gradient descent

$$s'(r) = \frac{d}{dr} \cdot \frac{1}{1+e^{-r}} = \frac{e^{-r}}{(1+e^{-r})^2} = s(r)(1-s(r))$$



$$\text{Let } S_i = s(x_i \cdot w)$$

$$\hookrightarrow \nabla_w J = - \sum_{i=1}^n \left(\frac{y_i}{S_i} \nabla S_i - \frac{1-y_i}{1-S_i} \cdot \nabla S_i \right)$$

$$= - \sum_{i=1}^n \left(\frac{y_i}{S_i} - \frac{1-y_i}{1-S_i} \right) S_i (1-S_i) x_i$$

$$= - \sum (y_i - S_i) x_i \quad \text{element-wise sigmoid of } x_i \text{ vector}$$

$$= - X^T (Y - S(Xw))$$

$$\text{gradient descent rule: } w \leftarrow w + \epsilon x^T (Y - S(Xw))$$

$$\text{SGD: } w \leftarrow w + \epsilon (y_i - s(x_i \cdot w)) x_i$$