

O. Big Picture

Statistical Learning IV
(26-11-2015)

①

Logarithmic loss:

* my predictions are probability distributions P for Y

$$* L(Y, P) = -\log P(Y) \in [0, \infty]$$

↳ small if I gave large probability to Y

↳ large if I gave small probability to Y .

Maximum Likelihood = ERM for log loss

$$\begin{aligned} \text{ML: } \arg \max_{\beta} \prod_{i=1}^N P_{\beta}(Y_i | X_i) &= \arg \min_{\beta} \sum_{i=1}^N -\log P_{\beta}(Y_i | X_i) \\ &= \arg \min_{\beta} \sum_{i=1}^N L(Y_i, P_{\beta}(\cdot | X_i)) \end{aligned}$$

↑
loss when learning is log loss
in big picture.

(1)

Statistical Learning IV (26-11-2015)

0. Big Picture
1. Surrogate losses
2. Logistic regression
3. Discriminative vs generative methods
- ~~unsupervised~~ Unsupervised learning
5. Clustering: K-means, EM with Gaussian mixtures

1. Surrogate losses

Classifications find \hat{f} with small expected 0/1-loss

$$EPE(\hat{f}) = \mathbb{E}_{x,y} [L(y, \hat{f}(x))]$$

$\begin{cases} 0 & \text{if } \hat{f}(x) \text{ predicts } y \text{ correctly} \\ 1 & \text{if } \hat{f}(x) \text{ predicts } y \text{ incorrectly} \end{cases}$

Empirical Risk Minimization:

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{arg\! min}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

Suppose \mathcal{F} set of all linear functions: $\{f(x) = x^\top \beta | \beta \in \mathbb{R}^p\}$

$$y \in \{-1, +1\} \quad L(y, f(x)) = \begin{cases} 0 & \text{if } \operatorname{sign}(f(x)) = y \\ 1 & \text{if } \operatorname{sign}(f(x)) \neq y \end{cases}$$

Difficulties:

- ERM solution not unique for linearly separable data
- computational problems for non-separable data

2

Surrogate/proxy loss L' :

Care about loss L , but train classifier with loss L'

Want: L large $\Leftrightarrow L'$ large

L small $\Leftrightarrow L'$ small

Usually: surrogate L' convex in parameters, so can optimize efficiently

Two most important surrogates: logistic loss, hinge loss

Predict correctly : iff $\text{sign}(x^T \beta) = y \Rightarrow y \cdot x^T \beta$ larger is better

(see Figure of surrogate losses)

2. Logistic Regression

(How can we use the power of linear regression for classification?)

Assume 2 classes for simplicity.

Bad idea: $\hat{p}(y=1|x) = x^T \hat{\beta}$ ← least squares for classification

E.g. y_i is going to the beach on day i .

x_{13} is temperature in °C

$$\hat{\beta}_3 = \frac{1}{5}$$

Then $+10^{\circ}\text{C} \Rightarrow x^T \beta + 2 \Rightarrow \hat{P}(y=1|x) + 2 > 1.0$

Probabilities do not behave linearly.

(3)

Better: $\log P_{\beta}(Y=1|X) = \mathbf{x}^T \beta \Leftrightarrow P_{\beta}(Y=1|X) = e^{\mathbf{x}^T \beta}$

$\mathbf{x}^T \beta + 2 \Rightarrow P_{\beta}(Y=1|X)$ increases by factor e^2 ,
but can still go outside $[0, 1]$

Logistic Regression:

$$\log \frac{P_{\beta}(Y=1|X)}{P_{\beta}(Y=-1|X)} = \mathbf{x}^T \beta \quad \text{log odds is linear} \quad \Leftrightarrow \begin{cases} P_{\beta}(Y=1|X) = \frac{e^{\mathbf{x}^T \beta}}{1+e^{\mathbf{x}^T \beta}} = \frac{1}{1+e^{-\mathbf{x}^T \beta}} \\ P_{\beta}(Y=-1|X) = \frac{1}{1+e^{\mathbf{x}^T \beta}} \\ P_{\beta}(Y|X) = \frac{1}{1+e^{-\mathbf{y}^T \beta}} \end{cases}$$

decision boundary:
x s.t. $\mathbf{x}^T \beta = 0$
is linear

logistic loss = log loss for logistic regression:

$$L(Y, \beta) = -\log P_{\beta}(Y|X) = \log(1+e^{-Y \mathbf{x}^T \beta})$$

Train with maximum likelihood = ERM for logistic loss

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^N \log(1+e^{-Y_i \mathbf{x}_i^T \beta})$$

No formula for minimum, but can find it by convex optimization

Important extensions:

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^N \log(1+e^{-Y_i \mathbf{x}_i^T \beta}) + \text{pen}(\beta)$$

$$L_2\text{-penalty: } \text{pen}(\beta) = \sum_{j=2}^p \beta_j^2 \quad (\text{like ridge regression})$$

$$L_1\text{-penalty: } \text{pen}(\beta) = \sum_{j=2}^p |\beta_j| \quad (\text{like lasso})$$

Considerations for penalization similar as for squared error.

Methods

④

3. Discriminative vs Generative Models

Probabilistic model: $\mathcal{P} = \{P_{\beta}(x, y) \mid \beta \in \mathbb{R}^P\}$

Generative

(imagine data generated by first choosing class y , then features x given y .)
e.g. Spam

estimate $P(x, y)$

LDA, Naive Bayes, ...

overkill for classification

Discriminative

estimate $P(y|x)$

Logistic regression, ...

less overkill, closer to class boundary, which is only thing we are interested in

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N -\log P_{\beta}(x_i, y_i)$$

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N -\log P_{\beta}(y_i | x_i)$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N -\log P_{\beta}(y_i | x_i) + \sum_{i=1}^N -\log P_{\beta}(x_i)$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N -\log P_{\beta}(y_i | x_i) + \sum_{i=1}^N -\log P^*(x_i)$$

acts as regularizer, so

reduces variance \Rightarrow better for small N

sample
does not depend
on $P(x)$

\rightarrow true distribution
of x

better fit of true distribution for large N
~~if $P^* \neq P$~~

Two methods are generative-discriminative pair if

one estimates $P_{\beta}(x, y)$ and other $P_{\beta}(y|x)$ in same model \mathcal{P} .

E.g.: ~~LDA~~ LDA - logistic regression for continuous features x_j
Naive Bayes - logistic regression for binary features

~~for discrete~~

(5)

Logistic Regression vs LDA

$x_{ij} \in \mathbb{R}$ (continuous features)

$$\text{LR: } \log \frac{\hat{P}(y=1|x)}{\hat{P}(y=-1|x)} = \mathbf{x}^T \hat{\beta} + \hat{\beta}_0 \quad \leftarrow \text{writing intercept separately}$$

$$\Leftrightarrow \hat{P}(y=1|x) = \frac{e^{\mathbf{x}^T \hat{\beta} + \hat{\beta}_0}}{1 + e^{\mathbf{x}^T \hat{\beta} + \hat{\beta}_0}}$$

$$\text{LDA: } \left\{ \begin{array}{l} \text{linear decision boundary} = \{x : \hat{P}(y=1|x) = \hat{P}(y=-1|x)\} \\ = \{x : \log \frac{\hat{P}(y=1|x)}{\hat{P}(y=0|x)} = 0\} \end{array} \right.$$

in fact, $\log \frac{\hat{P}(y=1|x)}{\hat{P}(y=-1|x)} = \mathbf{x}^T \hat{\alpha} + \hat{\alpha}_0$ is always linear in x ,
 ↗ not just for x s.t. it is 0
 parameters $\hat{\alpha}, \hat{\alpha}_0$ are
 a function of LDA parameters
 $\hat{\pi}_k, \hat{\Sigma}, \hat{\pi}_k$ ↗ possible to choose s.t. $\hat{\alpha} = \hat{\beta}$
 but LDA does not do that.
 $\hat{\alpha}_0 = \hat{\beta}_0$,
 so $\hat{P}(y=1|x) = \frac{e^{\mathbf{x}^T \hat{\alpha} + \hat{\alpha}_0}}{1 + e^{\mathbf{x}^T \hat{\alpha} + \hat{\alpha}_0}}$ ↗ has same form as
 for logistic regression.

LDA and LR are generative-discriminative pair

↳ expect LDA to be better for smaller N ,
 LR for larger N . (unless LDA model is correct)

Efron, 1975: if LDA model correct, then LR will perform
 as well as LDA if LR has 30% more data,
 so cannot be a lot worse.

(6)

Logistic Regression vs Naive Bayes

$x_{ij} \in \{0, 1\}$ (binary features), ~~so needs lots of features~~
 $j = 1, \dots, p$ ~~so very slow to calculate~~

\hookrightarrow transform into $x_i = \begin{pmatrix} n_0 \\ n_1 \end{pmatrix}$ where n_a is no. of a
 in i -th feature vector x_i

Will show that

$$(*) \quad \log \frac{\hat{P}(y=1|x)}{\hat{P}(y=-1|x)} = x^T \hat{\alpha} + \hat{\alpha}_0 \quad \text{for Naive Bayes (NB)}$$

determined by
 NB parameters
 for multinomial models $\hat{\alpha}_k$, $k=1, \dots, K$

Hence NB - LR form generative-discriminative pair.

So expect NB better for small N , LR better for large N .

NIPS

Ng, Jordan, 2001 make this precise

* Asymptotic performance as $N \rightarrow \infty$: LR always better than NB

* But ~~LR needs $N \geq O(p)$~~ to reach asymptotic performance

~~whereas~~ NB needs $N \geq O(\log p)$ to reach asymptotic performance.
 under some technical assumptions

$$\text{To show } (*): \log \frac{\hat{P}(y=1|x)}{\hat{P}(y=-1|x)} = \log \frac{\hat{P}(x|y=1) \hat{P}(y=1)}{\hat{P}(x|y=-1) \hat{P}(y=-1)} = \log \frac{\prod_i f_i(x_i) \pi_i}{\prod_i \tilde{f}_i(x_i) \tilde{\pi}_{i-1}}$$

f_k : multinomial model on binary features = Bernoulli model

simplify parametrization: $\hat{\alpha}_k$ is prob $x_{ij}=1$ for class k
 $-\hat{\alpha}_k$ is prob $x_{ij}=0$

$$\hat{P}(x|y=1) = \prod_i \hat{f}_i(x_i) = \prod_i \hat{\alpha}_k^{x_{ik}} (1 - \hat{\alpha}_k)^{1-x_{ik}}$$

$$+ \sum_{j=1}^p \hat{\alpha}_k^{x_{ij}} (1 - \hat{\alpha}_k)^{1-x_{ij}}$$

(7)

$$\begin{aligned}
 \log \frac{\hat{f}_k(x) \hat{\pi}_1}{\hat{f}_1(x) \hat{\pi}_1} &= \log \frac{\hat{\theta}_1^{n_1} (1-\hat{\theta}_1)^{n_0}}{\hat{\theta}_0^{n_0} (1-\hat{\theta}_0)^{n_1}} + \log \frac{\hat{\pi}_1}{\hat{\pi}_0} \\
 &= n_0 \cdot \log \frac{1-\hat{\theta}_1}{1-\hat{\theta}_0} + n_1 \cdot \log \frac{\hat{\theta}_1}{\hat{\theta}_0} + \log \frac{\hat{\pi}_1}{\hat{\pi}_0} \\
 &= x^T \hat{\alpha} + \hat{\alpha}_0 \quad \text{for } \hat{\alpha} = \begin{pmatrix} \log \frac{1-\hat{\theta}_1}{1-\hat{\theta}_0} \\ \log \frac{\hat{\theta}_1}{\hat{\theta}_0} \end{pmatrix} \quad \hat{\alpha}_0 = \log \frac{\hat{\pi}_1}{\hat{\pi}_0}
 \end{aligned}$$

possible to choose $\hat{\theta}_1, \hat{\theta}_0, \hat{\pi}_1, \hat{\pi}_0$, s.t. $\hat{\alpha} = \hat{\beta}$
 $\hat{\alpha}_0 = \hat{\beta}_0$,

but NB does not do that.

Unsupervised Learning

SS 43 Intro

Supervised: $T = (x_1, \dots, x_N)$ & can evaluate every method by EPE

Unsupervised: $T = x_1, \dots, x_N$

What can we do?

Clustering: split data into groups of points (clusters) that are similar

- ① What do you mean by "similar"?
- ② What kind of clusters are you looking for?

K-means

① Squared Euclidean distance between x_i, x_{i+1} :

$$d(x_i, x_{i+1}) = \sum_{j=1}^p (x_{ij} - x_{(i+1)j})^2$$

② Find K clusters s.t. distance to mean within each cluster is small.

$c(i) \in \{1, \dots, K\}$ is cluster assigned to x_i .

$$\text{minimize } \sum_{k=1}^K \sum_{i: c(i)=k} d(x_i, p_k) \quad (*)$$

where p_k = mean of points in cluster k

N.B. Mistake in book! (multiplies by N_k in each cluster)

Algorithm:

1. Initialize cluster assignment C
2. Set p_k to current mean in cluster k
3. Given means p_1, \dots, p_K , assign each x_i to nearest mean to get new C
4. repeat from 2 until no changes in C

- Converges to local minimum

- ~~often~~ ^{may} choose K s.t. increasing K does not reduce (*) very much.

Gaussian Mixtures and EM

- like k-means, but with "soft" cluster assignments

- for simplicity: explain for two clusters

cluster 1: $x \sim N(\mu_1, \sigma_1^2)$ w.p. $1-\pi$

cluster 2: $x \sim N(\mu_2, \sigma_2^2)$ w.p. π

probability density: $(1-\pi) \varphi_{\mu_1, \sigma_1^2}(x) + \pi \varphi_{\mu_2, \sigma_2^2}(x)$

each cluster
is a Gaussian
modelled
by

To find parameters:

ML is difficult numerically

+ gives bad solution with $\hat{\theta}_1 = 0$

$\hat{p}_i = x_i$ for some i .

Solution: algorithm very similar to K-means

Hidden variables $\Delta_i = \begin{cases} 0 & \text{if } x_i \text{ in cluster 1} \\ 1 & \text{if } x_i \text{ in cluster 2} \end{cases}$

1. Initializes parameters $\hat{\pi}, \hat{p}_1, \hat{\sigma}_1, \hat{p}_2, \hat{\sigma}_2$, ($\hat{\pi} > 0.5, \hat{\sigma}_i > 0.5$)

2. Soft cluster assignment:

$$\hat{\gamma}_i = \mathbb{E}[\Delta_i | \hat{\theta}, T] = \frac{\hat{\pi} \hat{q}_{\hat{p}_1, \hat{\sigma}_1}(x_i)}{\hat{\pi} \hat{q}_{\hat{p}_2, \hat{\sigma}_2}(x_i) + (1 - \hat{\pi}) \hat{q}_{\hat{p}_1, \hat{\sigma}_1}(x_i)}$$

3. Update parameter estimates for clusters:

$$\hat{p}_1 = \frac{\sum_{i=1}^N (\hat{\gamma}_i) x_i}{\sum_{i=1}^N (\hat{\gamma}_i)} \quad \hat{p}_2 = \frac{\sum_{i=1}^N \hat{\gamma}_i x_i}{\sum_{i=1}^N \hat{\gamma}_i}$$

$$\hat{\sigma}_1^2 = \frac{\sum_{i=1}^N (\hat{\gamma}_i)(x_i - \hat{p}_1)^2}{\sum_{i=1}^N (\hat{\gamma}_i)} \quad \hat{\sigma}_2^2 = \frac{\sum_{i=1}^N \hat{\gamma}_i(x_i - \hat{p}_2)^2}{\sum_{i=1}^N \hat{\gamma}_i}$$

$$\hat{\pi} = \frac{\sum_{i=1}^N \hat{\gamma}_i}{\sum_{i=1}^N \hat{\gamma}_i + \sum_{i=1}^N (1 - \hat{\gamma}_i)} = \frac{\sum_{i=1}^N \hat{\gamma}_i}{N}$$

4. repeat from 2 until convergence

- converges to local minimum.