

Statistical Learning IV

22-11-2019

1. Linear Discriminant Analysis (LDA)
2. Surrogate Losses
3. Log loss
4. Logistic Regression
5. Discriminative vs Generative Classification
 - ↳ Logistic Regression vs Naive Bayes

5. Linear Discriminant Analysis (LDA)

Model: $f_k(x)$ is $\mathcal{N}(\mu_k, \Sigma_k)$ is multivariate Gaussian
(see fig. 4.5)

Take $\Sigma_k = \Sigma$ the same for all k .

Parameter estimates from T :

$$\hat{\pi}_k = \frac{n_k}{N} \quad \text{where } n_k \text{ is number of observations in class } k.$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i: y_i=k} x_i \quad \text{is mean of } x_i \text{ in class } k$$

$$\hat{\Sigma} = \frac{1}{N-K} \sum_{k=1}^K \sum_{i: y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

$$\begin{aligned} \arg \max_k \hat{f}_k(x) \cdot \hat{\pi}_k &= \arg \max_k \log \hat{f}_k(x) + \log \hat{\pi}_k \\ &= \arg \max_k \log \left(\frac{1}{(2\pi)^{p/2} |\hat{\Sigma}|^{p/2}} e^{-\frac{1}{2}(x - \hat{\mu}_k)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_k)} \right) + \log \hat{\pi}_k \\ &= \arg \max_k -\frac{1}{2}(x - \hat{\mu}_k)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_k) + \log \hat{\pi}_k \\ &= \arg \max_k \underbrace{-\frac{1}{2} x^T \hat{\Sigma}^{-1} x + x^T \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k}_{\delta_k(x)} + \log \hat{\pi}_k \end{aligned}$$

linear in x ,

so decision boundary between any two classes also linear in x .

Hence name LDA

Decision boundary₀ (between classes k_1 and k_2)

$$x \text{ s.t. } \delta_{k_1}(x) - \delta_{k_2}(x) = 0$$

4. Surrogate Losses

Classification: find \hat{f} with small expected 0-1-loss

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

0 if $\hat{f}(x)$ predicts y correctly
1 if $\hat{f}(x)$ predicts y incorrectly

Empirical Risk Minimization:

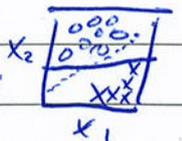
$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

Suppose \mathcal{F} is set of all linear functions: $\mathcal{F} = \{f_{\beta}(x) = x^T \beta \mid \beta \in \mathbb{R}^p\}$

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

Difficulties:

- ERM solution not unique for linearly separable data:



- Computational problems for non-separable data.

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

↑
today

↑
used in support vector machines

(next week)

For 0/1-loss, predict correctly iff

$$\text{sign}(x^T \hat{\beta}) = y \Leftrightarrow y x^T \hat{\beta} > 0$$

Most surrogate losses are therefore decreasing as a function of $y x^T \hat{\beta}$.

(see slide)

Exception: squared error

$$\text{N.B. } y^2 = 1$$

$$(y - x^T \hat{\beta})^2 = y^2 (y - x^T \hat{\beta})^2 = (y^2 - y x^T \hat{\beta})^2 = (1 - y x^T \hat{\beta})^2$$

increases for $y x^T \hat{\beta} > 1$.

This is undesirable!

5. Logarithmic Loss

* My predictions are probability distributions P for Y

$$* L(Y, P) = -\log P(Y)$$

↳ small if I gave large probability to Y

↳ large if I gave small probability to Y

Maximum Likelihood = ERM for log loss.

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

Bayes optimal predictor for log loss is true distribution:

$$\arg \min_P \mathbb{E}[L(Y, P)] = P^*(\cdot | X)$$

b. Logistic Regression (for two classes for simplicity)

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

Bad idea: $\hat{P}(Y=1 | X) = x^T \hat{\beta}$

E.g. Y_i is opening the window during i -th lecture

X_{i3} is avg temperature in the room in $^{\circ}\text{C}$.

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

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

Probabilities do not behave linearly.

(3)

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

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

Logistic Regression:

log odds is linear

$$\log \frac{P_{\beta}(Y=1|X)}{P_{\beta}(Y=-1|X)} = x^T \beta \Leftrightarrow \begin{cases} P_{\beta}(Y=1|X) = \frac{e^{x^T \beta}}{1 + e^{x^T \beta}} = \frac{1}{1 + e^{-x^T \beta}} \\ P_{\beta}(Y=-1|X) = \frac{1}{1 + e^{x^T \beta}} \end{cases}$$

$$P_{\beta}(Y|X) = \frac{1}{1 + e^{-y x^T \beta}}$$

decision boundary:

$$P_{\beta}(Y=1|X) > \frac{1}{2} \Leftrightarrow x^T \hat{\beta} > 0 \Rightarrow x \text{ s.t. } x^T \hat{\beta} = 0$$

is linear

logistic loss = log loss for logistic regression:

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

Train with maximum likelihood = ERM for logistic loss

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^N \log(1 + e^{-y_i x_i^T \beta})$$

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

Important extensions:

$$\text{minimize} \sum_{i=1}^N \log(1 + e^{-y_i x_i^T \beta}) + \lambda \text{pen}(\beta)$$

L_2 -penalty: $\text{pen}(\beta) = \sum_{j=1}^p \beta_j^2$ (like ridge regression)

L_1 -penalty: $\text{pen}(\beta) = \sum_{j=1}^p |\beta_j|$ (like lasso)

Considerations for penalization similar as for squared error.

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~~ for classification

Discriminative

estimate $P(y|x)$

Logistic regression, ...

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

$$\arg\min_{\beta} \sum_{i=1}^N -\log P_{\beta}(x_i, y_i)$$

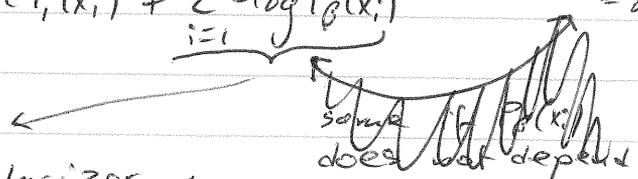
$$= \arg\min_{\beta} \sum_{i=1}^N -\log P_{\beta}(y_i | x_i) + \sum_{i=1}^N -\log P_{\beta}(x_i)$$

acts as regularizer, so reduces variance \Rightarrow better for small N

$$\arg\min_{\beta} \sum_{i=1}^N -\log P_{\beta}(y_i | x_i)$$

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

better fit of true distribution for large N if ~~model~~ $P^* \notin \mathcal{P}$



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

Logistic Regression vs Naive Bayes

$$\text{LR: } \log \frac{\hat{P}(y=1|X)}{\hat{P}(y=-1|X)} = X^T \hat{\beta} + \hat{\beta}_0$$

writing intercept separately

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

Will show that

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

determined by NB parameters $\hat{\theta}_{ij}, \hat{\pi}_j, \hat{\pi}_{-j}$

Hence Naive Bayes is working with a probabilistic model $\mathcal{P} = \{P_{\alpha}(X, y) | \alpha \in \mathbb{R}^p\}$ such that $P_{\alpha}(y|X)$ is the same as the model logistic regression is using.

So NB-LR form a generative-discriminative pair

Therefore expect NB better for smaller N ,
LR better for larger N .

Ng, Jordan, NIPS 2002 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 (under some technical assumptions) NB needs $N \geq O(\log p)$ to reach asymptotic performance

To show (*)₀

$$\log \frac{\hat{P}(y=1|x)}{\hat{P}(y=-1|x)} = \log \frac{\hat{P}(x|y=1) \hat{P}(y=1)}{\hat{P}(x|y=-1) \hat{P}(y=-1)} = \log \frac{f_1(x) \pi_1}{f_{-1}(x) \pi_{-1}}$$

$$x = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 2 \end{pmatrix} \begin{matrix} \leftarrow \text{UNIVERSITY} \\ \\ \leftarrow \text{BIG} \\ \\ \leftarrow \text{MONEY} \end{matrix} = \begin{pmatrix} n_1 \\ \vdots \\ n_p \end{pmatrix}$$

$$f_k(x) = \prod_{j=1}^p \theta_{kij}^{n_j}$$

$$\log \frac{f_1(x) \pi_1}{f_{-1}(x) \pi_{-1}} = \log \frac{\prod_{j=1}^p \theta_{1ij}^{n_j}}{\prod_{j=1}^p \theta_{-1ij}^{n_j}} + \log \frac{\pi_1}{\pi_{-1}}$$

$$= \log \prod_{j=1}^p \left(\frac{\theta_{1ij}}{\theta_{-1ij}} \right)^{n_j} + \log \frac{\pi_1}{\pi_{-1}}$$

$$= \sum_{j=1}^p n_j \log \frac{\theta_{1ij}}{\theta_{-1ij}} + \log \frac{\pi_1}{\pi_{-1}}$$

$$= x^T \tilde{\alpha} + \tilde{\alpha}_0 \quad \text{for } \tilde{\alpha} = \begin{pmatrix} \log \frac{\theta_{1,11}}{\theta_{-1,11}} \\ \vdots \\ \log \frac{\theta_{1,p}}{\theta_{-1,p}} \end{pmatrix}$$

$$\tilde{\alpha}_0 = \log \frac{\pi_1}{\pi_{-1}}$$

Possible to choose $\hat{\theta}_{kij}, \hat{\pi}_k$ such that $\tilde{\alpha} = \hat{\beta}, \tilde{\alpha}_0 = \hat{\beta}_0$, but NB does not do that \forall