Machine learning lecture slides

## Prediction theory

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Statistical model for binary outcomes
- Plug-in principle and IID model
- Maximum likelihood estimation
- Statistical model for binary classification
-
-
-

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Example: coin toss
- Physical model: hard
- Statistical model: outcome is random
  - *Bernoulli distribution* with heads probability $\theta \in [0,1]$



- Goal: correctly predict outcome

- Suppose $Y \sim \mathrm{Bernoulli}(\theta)$.
  - Suppose $\theta$ known.
  - Optimal prediction $\mathbf{1}_{\{\theta>1/2\}}$

  - The optimal prediction is incorrect with

- If $\theta$ unknown:
  - Assume we have data: outcomes of previous coin tosses
  - Data should be related to what we want to predict: same coin is being tossed

- *Plug-in principle*:
  - Estimate unknown(s) based on data (e.g., $\theta$)
  - Plug estimates into formula for optimal prediction

Assignment Project Exam Help

https://eduassistpro.github.i

  - *IID model*: Observations & (unseen) o_____ m variables
  - *iid* independent and identically distrib
  - Crucial modeling assumption that ma

Add WeChat edu_assist_pr

- When is the IID assumption not reasonable? . . .

## Statistical models

- *Parametric statistical model* $\{P_\theta : \theta \in \Theta\}$
  - collection of parameterized probability distributions for data
  - $\Theta$ is the *parameter space*
  - One distribution per parameter value $\theta \in \Theta$
- li

  https://eduassistpro.github.i (*pmf*)

  for the distribution.
  - What is formula for $P_\theta(y_1, \ldots, y\qquad\}^n$?

- *Likelihood* of parameter $\theta$ (given observed data)
  - $L(\theta) = P_\theta(y_1, \ldots, y_n)$
  - *Maximum likelihood estimation:*
  - Choose $\theta$ with highest likelihood
-   _____

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Coin toss example
  - Log-likelihood

$$\ln L(\theta) = \sum y_i \ln \theta + (1 - y_i) \ln(1 - \theta)$$

$$\hat{\theta}_{\text{MLE}} := -$$

- We are given data $y_1, \ldots, y_n \in \{0, 1\}^n$, which we model using the IID model from before
- Obtain estimate $\hat{\theta}_{\mathrm{MLE}}$ of unknown $\theta$ based on $y_1, \ldots, y_n$
- Plug-in $\theta_{\mathrm{MLE}}$ for $\theta$ in formula for optimal prediction:

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- How good is the plug-in prediction?
    - Study behavior under the IID model, where
    $Y_1, \ldots, Y_n, Y \sim_{\text{iid}} \text{Bernoulli}(\theta)$
        - $Y_1, \ldots, Y_n$ are the data we collected
        - $Y$ is the outcome to predict

https://eduassistpro.github.i

$Y$

worse.

Add WeChat edu_assist_pr

- **Theorem**:
  $$\Pr(\hat{Y} \neq Y) \leq \min\{\theta, 1-\theta\} + \frac{1}{2} \cdot |\theta - 0.5| \cdot e^{-2n(\theta-0.5)^2}.$$
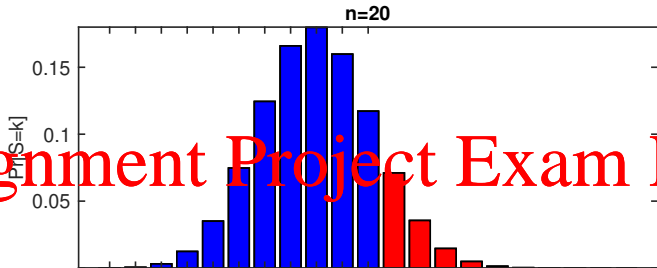  - The first term is the optimal error probability.
  - The second term comes from the probability that the $\hat{\theta}_{\mathrm{MLE}}$ is on the opposite side of $1/2$ as $\theta$.

Figure 1: $\Pr(S > n/2)$ for $S \sim$

Figure 2: $\Pr(S > n/2)$ for $S \sim$

Figure 3: $\Pr(S > n/2)$ for $S \sim$

Figure 4: $\Pr(S > n/2)$ for $S \sim$

- Example: spam filtering
- Labeled example: $(x, y) \in \mathcal{X} \times \{0, 1\}$
  - $\mathcal{X}$ is input (feature) space, $\{0, 1\}$ is the output (label) space
  - $\mathcal{X}$ is not necessarily the space of inputs itself (e.g., space of all

- $X$ has some *marginal probability d*
- *Conditional probability distribution* Bernoulli with heads probability
- $\eta \colon \mathcal{X} \to [0, 1]$ is a function, sometim *regression function* or *conditional mean function* (since $\overline{\mathbb{E}[Y \mid X = x] = \eta(x)}$).

- For a classifier $f\colon \mathcal{X} \to \{0,1\}$, the *error rate* of $f$ (with respect to the distribution of $(X,Y)$) is

$$\mathrm{err}(f) := \Pr(f(X) \neq Y),$$

$$\frac{1}{| \ |} \sum_{(x,y)} \qquad$$

which is the same as $\Pr(f(X) \neq Y)$
$(X,Y)$ is uniform over the labeled exampl

- Caution: This notation $\mathrm{err}(f)$ does not make explicit the dependence on (the distribution of) the random example $(X,Y)$. You will need to determine this from context.

- Consider any random variables $A$ and $B$.
- Conditional expectation of $A$ given $B$:
  - Written $\mathbb{E}[A \mid B]$
  - A random variable! What is its expectation?
  - *Law of iterated expectations* (a.k.a. *tower property*):

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Example: roll a fair 6-sided die
  - $A =$ number shown facing up
  - $B =$ parity of number shown facing up
  - $C := \mathbb{E}[A \mid B]$ is random variable with

$$= \frac{1}{2}$$

$$= \frac{1}{2}$$

- *Optimal classifier* (*Bayes classifier*):

$$f^\star(x) = \mathbf{1}_{\{\eta(x) \geq 1/2\}}$$

where $\eta$ is the conditional mean function

- ▶ Write error rate as $\mathrm{err}(f^\star) = \mathrm{Pr}$ ${}^\star$ / $_{/Y\}}]$
  - ▶ Conditional on $X$, probability of mis
  - ▶ $\min\{\eta(X), 1-\eta(X)\}$
  - ▶ So, optimal error rate is

$$\mathrm{err}(f^\star) = \mathbb{E}[\mathbf{1}_{\{f^\star(X) \neq Y\}}]$$
$$= \mathbb{E}[\mathbb{E}[\mathbf{1}_{\{f^\star(X) \neq Y\}} \mid X]]$$
$$= \mathbb{E}[\min\{\eta(X), 1-\eta(X)\}].$$

- Suppose input $x$ is a single (binary) feature, "is email all-caps?"
- How to interpret "the probability that email is spam given

Assignment Project Exam Help

- What does it mean for the Bayes classifier $f^*$ to be optimal?

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- What to do if $\eta$ is unknown?
  - Training data: $(x_1, y_1), \ldots, (x_n, y_n)$
  - Assume data are related to what we want to predict
  - Let $Z := (X, Y)$, and $Z_i := (X_i, Y_i)$ for $i = 1, \ldots, n$.
  - IID model: $Z_1, \ldots, Z_n, Z$ are iid random variables

- Study in context of IID model
- Assume $\eta(x) \approx \eta(x')$ whenever $x$ and $x'$ are close.
  - This is where the modeling assumption comes in (via choice of distance function).
- Let $(X, Y)$ be the "test" example, and suppose $(X_{\hat{i}}, Y_{\hat{i}})$ is the

  https://eduassistpro.github.i

  $\eta(X) \approx \eta(X_{\hat{i}})$.
- Prediction is $Y_{\hat{i}}$, true label is $Y$.
- Conditional on $X$ and $X_{\hat{i}}$, what is pro Add WeChat edu_assist_pr
  - $\eta(X)(1 - \eta(X_{\hat{i}})) + (1 - \eta(X))$
- Conclusion: expected error rate is
  $\mathbb{E}[\mathrm{err}(\mathrm{NN}_S)] \approx 2 \cdot \mathbb{E}[\eta(X)(1 - \eta(X))]$ for large $n$
  - Recall that optimal is $\mathbb{E}[\min\{\eta(X), 1 - \eta(X)\}]$.
  - So $\mathbb{E}[\mathrm{err}(\mathrm{NN}_S)]$ is at most twice optimal.
  - Never exactly optimal unless $\eta(x) \in \{0, 1\}$ for all $x$.

- How to estimate error rate?
- IID model:
  - $(X_1, Y_1), \ldots, (X_n, Y_n), (X'_1, Y'_1), \ldots, (X'_m, Y'_m), (X, Y)$ are iid.
  - Training examples (that you have): $(X_1, Y_1), \ldots, (X_n, Y_n)$
  - [Test examples: $(X'_1, Y'_1), \ldots, (X'_m, Y'_m)$]
- [Hence, **test examples are independent** of training examples (this is important!)]
- We would like to estimate $\mathrm{err}(\hat{f})$
  - Caution: since $\hat{f}$ depends on training
  - Convention: When we write $\mathrm{err}(\hat{f})$ where $\hat{f}$ is random, we really mean $\Pr(\hat{f}(X) \neq Y \mid \hat{f})$.
  - Therefore $\mathrm{err}(\hat{f})$ is a random variable!

- Conditional distribution of $S := \sum_{i=1}^m \mathbf{1}_{\{\hat{f}(X'_i) \neq Y'_i\}}$ given training data:

- $S \mid$ training data $\sim \text{Binomial}(m, \varepsilon)$ where $\varepsilon = \text{err}(\hat{f})$
- By law of large numbers,

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{}$$

  is close to $\varepsilon$ when $m$ is large
- How accurate is the estimate? Depends on the (conditional) variance!
  - $\text{var}(\frac{1}{m} S \mid \text{training data}) = \frac{\varepsilon(1-\varepsilon)}{m}$
  - Standard deviation is $\sqrt{\frac{\varepsilon(1-\varepsilon)}{m}}$

- *True positive rate* (*recall*): $\Pr(f(X) = 1 \mid Y = 1)$
- *False positive rate*: $\Pr(f(X) = 1 \mid Y = 0)$
- *Precision*: $\Pr(Y = 1 \mid f(X) = 1)$
- ...
- _____

Assignment Project Exam Help

https://eduassistpro.github.i

| | | |
|---|---|---|
| $y = 1$ | # false negatives | |

Add WeChat edu_assist_pr

- *Receiver operating characteristic (ROC) curve*
  - What points are achievable on the TPR-FPR plane?
  - Use randomization to combine classifiers
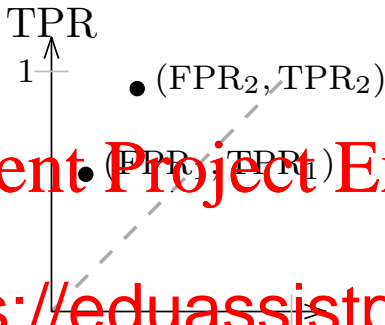
Assignment Project Exam Help

https://eduassistpro.github.i
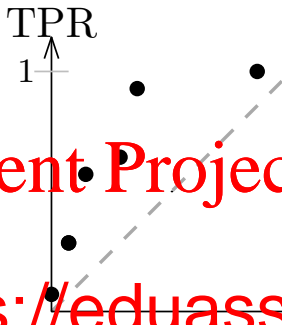
Add WeChat edu_assist_pr

Figure 5: TPR vs FPR plot with two poi

Figure 6: TPR vs FPR plot with many po

- What if there are $K > 2$ possible outcomes?
- Replace coin with $K$-sided die
- Say $Y$ has a *categorical distribution* over $[K] = \{1, \ldots, K\}$ determined probability vector $\theta = (\theta_1, \ldots, \theta_K)$

- https://eduassistpro.github.i

$$\hat{y} := \underset{k \in [}{\arg \max}$$

Add WeChat edu_assist_pr

- Statistical model for labeled examples $(X, Y)$, where $Y$ takes values in $[K]$
  - Now $Y \mid X = x$ has a categorical distribution with parameter vector $\eta(x) = (\eta(x)_1, \ldots, \eta(x)_K)$
  - Conditional probability function: $\eta(x)_k := \Pr(Y = k \mid X = x)$

- Example: Train OCR digit classifier using data from Alice's handwriting, but eventually use on digits written by Bob.
  - What is a better evaluation?

- What if we want to eventually use on digits wri Alice and Bob?