# 4.1 Evaluation Metrics

Qingqing Huang, Mu Li, Alex Smola

https://c.d2l.ai/stanford-cs329p

# Model Metrics

- Loss measures how good the model in predicting the outcome in supervised learning

- Other metrics to evaluate the model performance

  - Model specific: e.g. accuracy for classification, mAP for object detection

  - Business specific: e.g. revenue, inference latency

- We select models by multiple metrics

  - Just like how you choose cars



Small Cars    Midsize Cars    Large Cars    SUVs    Crossovers    Trucks    Hybrids    Vans    CPO    Luxury    Sports Cars    Convertibles

# Metrics for Binary <mark>Classification</mark>

- Accuracy: # correct predictions / # examples

```
sum(y == y_hat) / y.size
```

- Precision: # True positive / # (True positive + False positive)

```
sum((y_hat == 1) & (y == 1)) / sum(y_hat == 1)
```

- Recall: # True positive / # Positive examples
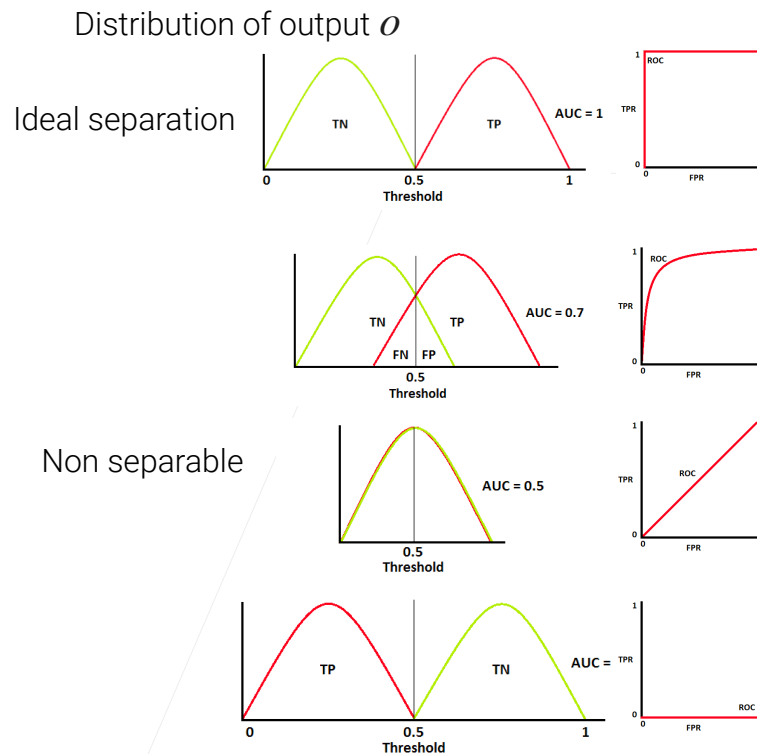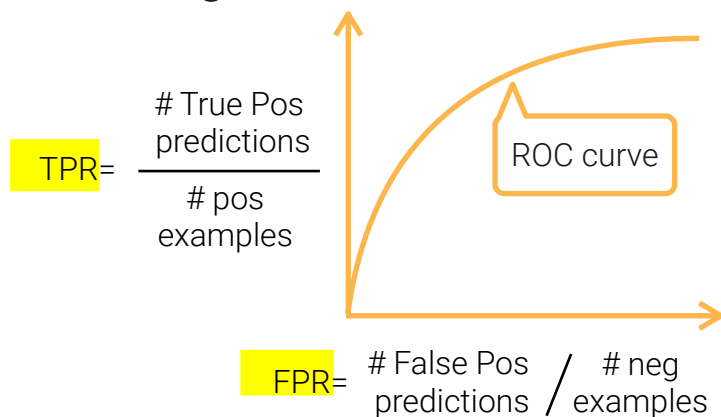
```
sum((y_hat == 1) & (y == 1)) / sum(y == 1)
```

- Be careful of <mark>division by 0</mark>

- One metric that balances precision and recall

  - F1: the harmonic mean of precision and recall: <mark>$2pr/(p + r)$</mark>
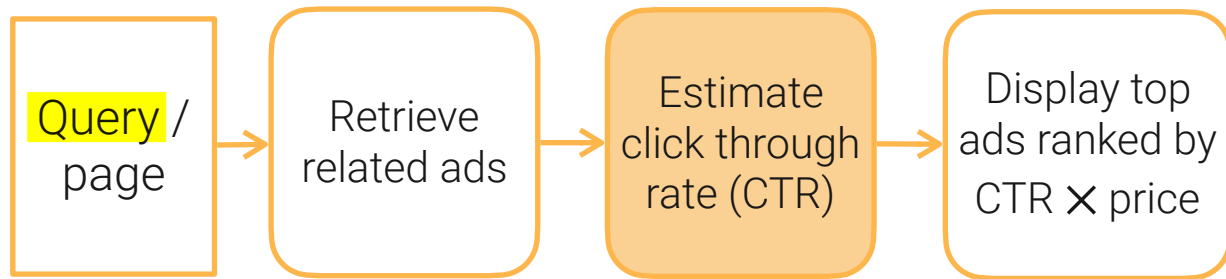
# AUC-ROC

- Measures how well the model can separate the two classes

- Choose decision threshold $\theta$, predict positive if $o \geq \theta$ else neg

- In the range [0.5, 1]

$$TPR = \frac{\text{\# True Pos predictions}}{\text{\# pos examples}}$$

ROC curve

$$FPR = \frac{\text{\# False Pos predictions}}{\text{\# neg examples}}$$

Distribution of output $o$

Ideal separation

Non separable

# Case Study: Displaying Ads

- Ads is one major revenue source for Internet companies



```
Query / page  →  Retrieve related ads  →  Estimate click through rate (CTR)  →  Display top ads ranked by CTR ✕ price  →
```

A binary classification problem

# Business Metrics for Displaying Ads

- Optimize both revenue and customer experience
  - Latency: ads should be shown to users at the same time as others
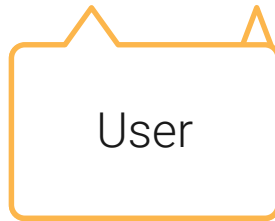  - ASN: average #ads shown in a page
  - CTR: actual user click through rate
  - ACP: average price advertiser pays per click
- revenue = #pageviews × ASN × CTR × ACP

Platform company   Matters to whom   User   Advertiser

# Displaying Ads: Model → Business Metrics

- The key model metric is AUC

- A new model with increased AUC may harm business metrics, possible reasons:

  - Lower estimated CTR → less ads displayed
  - Lower real CTR because we trained and evaluated on past data
  - Lower prices

- Online experiment: deploy models to evaluate on real traffic data

# Summary

- We evaluate models with ==multiple metrics==

- Model metrics evaluate model performance on ==examples==

    - E.g. ==accuracy, precision, recall, F1, AUC== for classification models

- ==Business metrics== measure how models impact the ==product==