CS-GY 6923: Lecture 3
Model Selection + Regularization + Bayesian
Perspective

NYU Tandon School of Engineering, Prof. Christopher Musco

- Homework 1 due tonight.
- New lab will be released tonight, due next Thursday.
- Next problem set will be due a week after that.

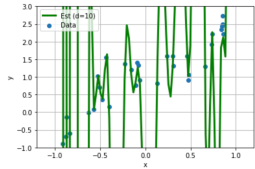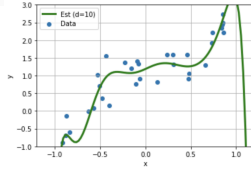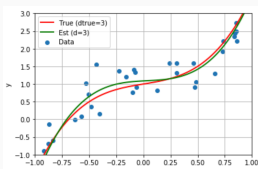Basic machine learning problem:

- Given model $f_{\boldsymbol{\theta}}$ and <u>loss function</u> $L(f_{\boldsymbol{\theta}})$.
- Choose $\boldsymbol{\theta}^*$ to minimize $L(f_{\boldsymbol{\theta}})$.

Model selection problem:

- Given choice of many models $f^{(1)}_{\boldsymbol{\theta}_1}, f^{(2)}_{\boldsymbol{\theta}_2}, \ldots, f^{(q)}_{\boldsymbol{\theta}_q}$.
- Choose $\boldsymbol{\theta}^*_1, \ldots, \boldsymbol{\theta}^*_q$ to minimize $L(f_{\boldsymbol{\theta}_1}), \ldots, L(f_{\boldsymbol{\theta}_q})$.
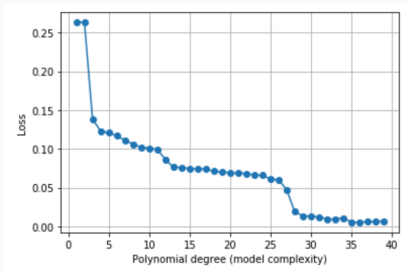- Then choose the "best" model for our data.

Polynomial regression models with different degree. See
`demo_polyfit.ipynb`.



- Model $f_{\boldsymbol{\theta}_1}^{(1)}$: all linear functions.
- Model $f_{\boldsymbol{\theta}_2}^{(2)}$: all quadratic functions.
- Model $f_{\boldsymbol{\theta}_3}^{(3)}$: all cubic functions.
- ...

The more **complex** our model class (e.g., the higher degree we allow in polynomial regression) the better our loss:
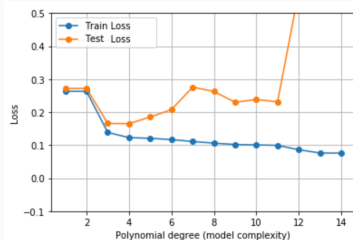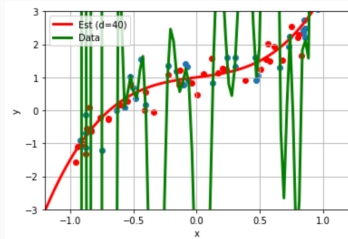


Training loss alone is not usually a good metric for model selection. Small loss does not imply generalization to new data.

Main approach: Evaluate model on fresh <u>test data</u> which was not used during training.

Test/train split:

- Given data set $(X, y)$, split into two sets $(X_{train}, y_{train})$ and $(X_{test}, y_{test})$.
- Train $q$ models $f^{(1)}, \ldots, f^{(q)}$ by finding parameters which minimize the loss on $(X_{train}, y_{train})$.
- Evaluate loss of each trained model on $(X_{test}, y_{test})$.

While train error always decreases, we eventually see test error
increase with increasing model complexity.

The above trend is fairly representative of what we tend to see across the board:

Is "test error" the end goal though? Don't we care about "future" error?

Intuition: Models which perform better on the test set will generalize better to future data.

Goal: Introduce a little bit of formalism to better understand what this means. What is "future" data?

**Statistical Learning Model:**

- Assume each data example is randomly drawn from some distribution $(\mathbf{x}, y) \sim \mathcal{D}$.



E.g. $x_1, \ldots, x_d$ are Gaussian random variables with parameters $\mu_1, \sigma_1, \ldots, \mu_d, \sigma_d$.

This is not a simplifying assumptions! The distribution could be arbitrarily complicated.

Statistical Learning Model:

- Assume each data example is randomly drawn from some distribution $(\mathbf{x}, y) \sim \mathcal{D}$.

- Define the **Risk** of a model/parameters:

$$R(f, \boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ L\left( f(\mathbf{x}, \boldsymbol{\theta}), y \right) \right]$$

here $L$ is our loss function (e.g. $L(z, y) = |z - y|$ or $L(z, y) = (z - y)^2$).

**Goal:** Find model $f \in \{f^{(1)}, \ldots, f^{(q)}\}$ and parameter vector $\boldsymbol{\theta}$ to minimize the $R(f, \boldsymbol{\theta})$.

- (Population) Risk:

$$R(f, \boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \left[ L\left( f(\mathbf{x}, \boldsymbol{\theta}), y \right) \right]$$

- Empirical Risk: Draw $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n) \sim \mathcal{D}$

$$R_E(f, \boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} L\left( f(\mathbf{x}, \boldsymbol{\theta}), y \right)$$

For any <u>fixed</u> model $f$ and parameters $\boldsymbol{\theta}$,

$$\mathbb{E}\left[R_E(f, \boldsymbol{\theta})\right] = R(f, \boldsymbol{\theta}).$$

Only true if $f$ and $\boldsymbol{\theta}$ are chosen *without looking at the data used to compute the empirical risk.*

- Train $q$ models $(f^{(1)}, \boldsymbol{\theta}_1^*), \ldots, (f^{(q)}, \boldsymbol{\theta}_q^*)$.
- For each model, compute empirical risk $R_E(f^{(i)}, \boldsymbol{\theta}_i^*)$ using <u>test data</u>.
- Since we assume our original dataset was drawn independently from $\mathcal{D}$, so is the random test subset.

No matter how our models were trained or how complex they are, $R_E(f^{(i)}, \boldsymbol{\theta}_i^*)$ is an <u>unbiased estimate</u> of the true risk $R(f^{(i)}, \boldsymbol{\theta}_i^*)$ for every $i$. Can use it to distinguish between models.

**bag-of-words** models and **n-grams**

Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.

This is a sentence.

a aardvark                                                              zoo zyzzyva

**bag-of-words**

**bag-of-words** models and **n-grams**

Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.



This is a sentence.

| | 1 | | | | | | | | 1 | | | | | | | | | | | | | | | 1 | | | | | | | | |

"a sentence"          "is a"                                          "this is"

## bi-grams

**bag-of-words** models and **n-grams**

Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.
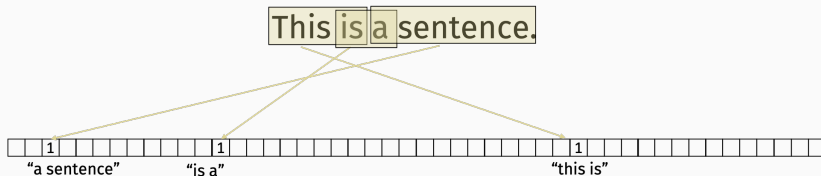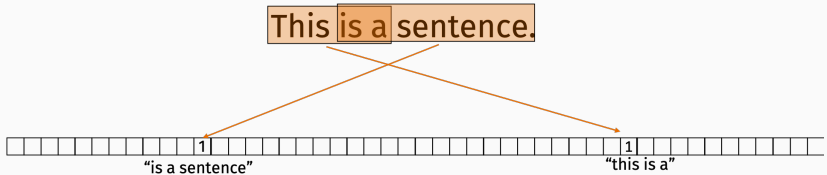


tri-grams

Models of increasing order:

- Model $f_{\boldsymbol{\theta}_1}^{(1)}$: spam filter that looks at **single words**.
- Model $f_{\boldsymbol{\theta}_2}^{(2)}$: spam filter that looks at **bi-grams**.
- Model $f_{\boldsymbol{\theta}_3}^{(3)}$: spam filter that looks at **tri-grams**.
- . . .
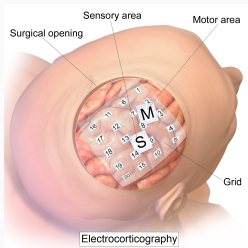
  "interest"          "low interest"          "low interest loan"

Increased length of **n-gram** means more expressive power.
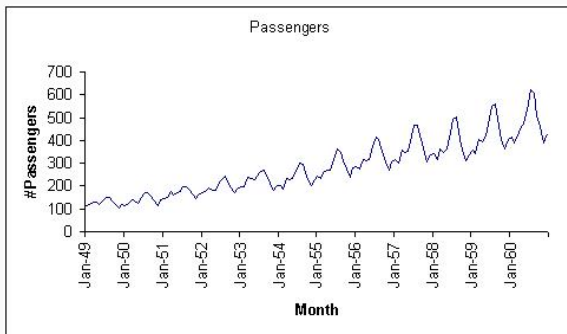
**Electrocorticography ECoG (upcoming lab):**

- Implant grid of electrodes on surface of the brain to measure electrical activity in different regions.



- Predict hand motion based on ECoG measurements.
- **Model order:** predict movement at time $t$ using brain signals at time $t, t-1, \ldots, t-q$ for varying values of $q$.

Predicting time $t$ based on a linear function of the signals at time $t, t-1, \ldots, t-q$ is <u>not the same</u> as fitting a line to the time series. It's much more expressive.



Passengers

Predecessor of modern "recurrent neural networks".

Electrocorticography ECoG lab:



**First lab where computation actually matters (solving regression problems with $\sim 40k$ examples, $\sim 1500$ features)**

Makes sense to test and debug code using a subset of the data.

Slight caveat: This is typically not how machine learning or scientific discovery works in practice!

Typical workflow:

- Train a class of models.
- Test.
- Adjust class of models.
- Test.
- Adjust class of models.
- Cont...

Final model implicitly depends on test set because performance on the test set guided how we changed our model.

# Popularity of ML benchmarks and competitions leads to adaptivity at a massive scale.



Kaggle (various competitions)



Imagenet (image classification and categorization)

## Is adaptivity a problem? Does it lead to over-fitting? How much? How can we prevent it? All current research.

# The reusable holdout: Preserving validity in adaptive data analysis

Cynthia Dwork[1,*], Vitaly Feldman[2,*], Moritz Hardt[3,*], Toniann Pitassi[4,*], Omer Reingold[5,*], Aaron Roth[6,*]
+ See all authors and affiliations

## Do ImageNet Classifiers Generalize to ImageNet?

Benjamin Recht*    Rebecca Roelofs    Ludwig Schmidt    Vaishaal Shankar
UC Berkeley         UC Berkeley        UC Berkeley       UC Berkeley

12 Jun 2019

**Abstract**

We build new test sets for the CIFAR-10 and ImageNet datasets. Both benchmarks have been the focus of intense research for almost a decade, raising the danger of overfitting to excessively re-used test sets. By closely following the original dataset creation processes, we test to what extent current classification models generalize to new data. We evaluate a broad range of models and find accuracy drops of $3\% - 15\%$ on CIFAR-10 and $11\% - 14\%$ on ImageNet. However, accuracy gains on the original test sets translate to larger gains on the new test sets. Our results suggest that the accuracy drops are not caused by adaptivity, but by the models' inability to generalize to slightly "harder" images than those found in the original test sets.

25

## Do ImageNet Classifiers Generalized to ImageNet?



Interestingly, when comparing popular vision models on "fresh" data, while performance dropped across the board, the relative rank of model performance did not change significantly.

REGULARIZATION

In all the model selection examples we discussed we had full control over the complexity of the model: could range from underfitting to overfitting.

In practice, you often don't have this freedom. Even the most basic model might lead to overfitting.



d features

n examples

$X_1$
$X_2$
⋮
$X_n$

$X$

$y_1$
$y_2$
⋮
$y_n$

$y$

**Example:** Linear regression model where $d \geq n$. Can always find $\beta$ so that $X\beta = y$ exactly.

Select some subset of features to use in model:



Filter method: Compute some metric for each feature, and select features with highest score.

- Example: compute loss or $R^2$ value when each feature in X is used in single variate regression.

Any potential limitations of this approach?

**Exhaustive approach:** Pick best subset of $q$ features.

**Faster approach:** Greedily select $q$ features.

Stepwise Regression:

- **Forward:** Step 1: pick single feature that gives lowest loss. Step $k$: pick feature that when combined with previous $k - 1$ chosen features gives lowest loss.
- **Backward:** Start with all of the features. Greedily eliminate those which have least impact on model performance.

Feature selection deserves more than two slides, but we won't go into too much more detail!

**Regularization:** Explicitly discourage overfitting by adding a regularization penalty to the loss minimization problem.

$$\min_{\boldsymbol{\theta}} \left[ L(\boldsymbol{\theta}) + Reg(\boldsymbol{\theta}) \right].$$

**Example:** Least squares regression. $L(\boldsymbol{\beta}) = \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2$.

- Ridge regression ($\ell_2$): $Reg(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_2^2$
- LASSO (least absolute shrinkage and selection operator) ($\ell_1$): $Reg(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_1$
- Elastic net: $Reg(\boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2^2$

Ridge regression: $\min_{\boldsymbol{\beta}} \|X\boldsymbol{\beta} - y\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2$.

- As $\lambda \to \infty$, we expect $\|\boldsymbol{\beta}\|_2^2 \to 0$ and $\|X\boldsymbol{\beta} - y\|_2^2 \to \|y\|_2^2$.
- Feature selection methods attempt to set many coordinates in $\boldsymbol{\beta}$ to 0. Ridge regularizations encourages coordinates to be small.

Ridge regression: $\min_{\boldsymbol{\beta}} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2$.

· Can be viewed as shrinking the size of our model class. Relaxed version of $\min_{\boldsymbol{\beta}:\|\boldsymbol{\beta}\|_2^2 < c} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2$.

**Claim:** For any $\lambda$, let $\boldsymbol{\beta}_\lambda^* = \arg\min_{\boldsymbol{\beta}} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2$. Then there is some $c(\lambda)$ such that:

$$\boldsymbol{\beta}_\lambda^* = \underset{\boldsymbol{\beta}:\|\boldsymbol{\beta}\|_2^2 < c(\lambda)}{\arg\min} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2.$$

Moreover, we have the for $\lambda' > \lambda$, $c(\lambda') < c(\lambda)$.

Ridge regression: $\min_{\boldsymbol{\beta}} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2$.

- $\min_{\boldsymbol{\beta}:\|\boldsymbol{\beta}\|_2^2 < c} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2$ won't have a solution at zero for all $\mathbf{y}$, even when over-parameterized.



- Regularization methods are <u>not invariant</u> to data scaling. Typically when using regularization we mean center and scale columns to have unit variance.

How do we minimize: $L_R(\boldsymbol{\beta}) = \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2$?

Lasso regularization: $\min_{\boldsymbol{\beta}} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$.

- As $\lambda \to \infty$, we expect $\|\boldsymbol{\beta}\|_1 \to 0$ and $\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 \to \|\mathbf{y}\|_2^2$.
- Typically encourages subset of $\boldsymbol{\beta}_i$'s to go to zero, in contrast to ridge regularization.

Pros:

- Simpler, more interpretable model.
- More intuitive reduction in model order.

Cons:

- No closed form solution because $\|\boldsymbol{\beta}\|_1$ is not differentiable.
- Can be solved with iterative methods, but generally not as quickly as ridge regression.

Notes:

- Model selection/cross validation used to choose optimal scaling $\lambda$ on $\lambda\|\boldsymbol{\beta}\|_2^2$ or $\lambda\|\boldsymbol{\beta}\|_1$.
- Often grid search for best parameters is performed in "log space". E.g. consider $[\lambda_1, \ldots, \lambda_q] = 1.5^{[-4,-3,-2,-1,-0,1,2,3,4]}$.

# THE BAYESIAN/PROBABILISTIC MODELING PERSPECTIVE

- Data Examples: $x_1, \ldots, x_n \in \mathbb{R}^d$
- Target: $y_1, \ldots, y_n \in \{0, 2, \ldots, q-1\}$ when there are $q$ classes.
    - Binary Classification: $q = 2$, so each $y_i \in \{0, 1\}$.
    - Multi-class Classification: $q > 2$. [1]

---

[1]Note that there is also <u>multi-label</u> classification where each data example maybe belong to more than one class.

- Medical diagnosis from MRI: 2 classes.
- MNIST digits: 10 classes.
- Full Optical Character Regonition: 100s of classes.
- ImageNet challenge: 21,000 classes.

Running example today: Email Spam Classification.

Classification can (and often is) solved using the same **loss-minimization framework** we saw for regression.

We won't see that today! We're going to use classification as a window into another way of thinking about machine learning.

Will give new an interesting justifications for tools like <u>regularization.</u>

**Today:** ML from a Probabilistic Modeling/Bayesian Perspective.

In a <u>Bayesian</u> or <u>Probabilistic</u> approach to machine learning we always start by conjecturing a

<div align="center">probabilistic model</div>

that plausibly could have generated our data.

- The model guides how we make predictions.
- The model typically has unknown parameters $\vec{\theta}$ and we try to find the most reasonable parameters based on observed data (more on this later in lecture).

Typically we try to keep things simple!

**Exercise:** Come up with a probabilistic model for <u>any one</u> of the following data sets $(x_1, y_1), \ldots, (x_n, y_n)$.

1. For $n$ **people**: each $x_i \in \{0, 1\}$ with zero indicating <u>male</u>, one indicating <u>female</u>. Each $y_i$ is the height of the person in inches.
2. For $n$ **NYC apartments**: each $x_i$ is the size of the apartment in square feet. Each $y_i$ is the monthly rent in dollars.
3. For $n$ **students**: each $x_i \in \{Fresh., Soph., Jun., Sen.\}$ indicating class year. Each $y_i \in \{0, 1\}$ with zero indicating the student has not taken machine learning, one indicating they have.

What are the unknown parameters of your model. What would be a guess for their values? How would you confirm or refine this guess using data?

Dataset: $(x_1, y_1), \ldots, (x_n, y_n)$

Description: For *n* people: each $x_i \in \{0, 1\}$ with zero indicating
<u>male</u>, one indicating <u>female</u>. Each $y_i$ is the height of the person
in inches.

Model:

Dataset: $(x_1, y_1), \ldots, (x_n, y_n)$

Description: For $n$ NYC apartments: each $x_i$ is the size of the apartment in square feet. Each $y_i$ is the monthly rent in dollars.

Model:

Dataset: $(x_1, y_1), \ldots, (x_n, y_n)$

Description: For $n$ students: each
$x_i \in \{Fresh., Soph., Jun., Sen.\}$ indicating class year. Each
$y_1 \in \{0, 1\}$ with zero indicating the student has not taken
machine learning, one indicating they have.

Model:

Goal:

- Build a probabilistic model for a binary classification problem.
- Estimate parameters of the model.
- From the model derive a classification rule for future predictions (the Naive Bayes Classifier).

feature extraction

bag-of-words

ML prediction

1 0 1 1 1 0 0 0 0 0 0 1 1 1 0 → 0 (safe)

1 0 0 0 1 0 1 0 0 1 0 0 0 0 0 → 1 (spam)

1 0 0 0 0 1 1 1 0 0 0 0 0 0 0 → 0 (safe)

1 0 0 0 0 1 0 0 1 0 1 0 0 0 0 → 0 (safe)

1 0 0 0 1 0 1 0 1 0 0 1 1 0 0 → 1 (spam)

Both target labels and data vectors are binary.

Probabilistic model for (bag-of-words, label) pair $(\mathbf{x}, y)$:

- Set $y = 0$ with probability $p_0$, $y = 1$ with probability $p_1 = 1 - p_0$.
  - $p_0$ is probability an email is not spam (e.g. 99%).
  - $p_1$ is probability an email is spam (e.g. 1%).
- If $y = 0$, for each $i$, set $x_i = 1$ with prob. $p_{i0}$.
- If $y = 1$, for each $i$, set $x_i = 1$ with prob. $p_{i1}$.

Unknown model parameters:

- $p_0, p_1$,
- $p_{10}, p_{20}, \ldots p_{n0}$, one for each of the $n$ vocabulary words.
- $p_{11}, p_{21}, \ldots p_{n1}$, one for each of the $n$ vocabulary words.

How would you estimate these parameters?

Reasonable way to set parameters:

- Set $p_0$ and $p_1$ to the empirical fraction of not spam/spam emails.
- For each word $i$, set $p_{i0}$ to the empirical probability word $i$ appears in a <u>non-spam</u> email.
- For each word $i$, set $p_{i1}$ to the empirical probability word $i$ appears in a <u>spam</u> email.

Estimating these parameters from previous data examples is the only "training" we will do.

# DONE WITH MODELING
## ON TO PREDICTION

- **Probability:** p(x) – the probability event *x* happens.
- **Joint probability:** p(x,y) – the probability that event *x* <u>and</u> event *y* happen.
- **Conditional Probability** *p(x | y)* – the probability *x* happens <u>given</u> that *y* happens.

$$p(x|y) =$$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Proof:

Given unlabeled input $(\mathbf{x}, \underline{\quad})$, choose the label $y \in \{0, 1\}$ which is <u>most likely</u> given the data. Recall $\mathbf{x} = [0, 0, 1, \ldots, 1, 0]$.

**Classification rule:** maximum a posterior prob. (MAP) estimate.

Step 1. Compute:

- $p(y = 0 \mid \mathbf{x})$: prob. $y = 0$ given observed data vector $\mathbf{x}$.
- $p(y = 1 \mid \mathbf{x})$: prob. $y = 1$ given observed data vector $\mathbf{x}$.

**Step 2. Output:** 0 or 1 depending on which probability is larger.

$p(y = 0 \mid \mathbf{x})$ and $p(y = 1 \mid \mathbf{x})$ are called **posterior** probabilities.

How to compute the posterior? **Bayes rule!**

$$p(y = 0 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 0)p(y = 0)}{p(\mathbf{x})} \tag{1}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \tag{2}$$

- **Prior:** Probability in class 0 <u>prior</u> to seeing any data.
- **Posterior:** Probability in class 0 <u>after</u> seeing the data.

Goal is to determine which is larger:

$$p(y = 0 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 0)p(y = 0)}{p(\mathbf{x})} \qquad \text{vs.}$$

$$p(y = 1 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 1)p(y = 1)}{p(\mathbf{x})}$$

How to compute posteriors:

- Ignore evidence $p(\mathbf{x})$ since it is the same for both sides.
- $p(y = 0)$ and $p(y = 1)$ already known (computed from training data).
- $p(\mathbf{x} \mid y = 0) = ?$ $p(\mathbf{x} \mid y = 1) = ?$

"Naive" Bayes Rule: Compute $p(\mathbf{x} \mid y = 0)$ by assuming
<u>independence</u>:

$$p(\mathbf{x} \mid y = 0) = p(x_1 \mid y = 0) \cdot p(x_2 \mid y = 0) \cdot \ldots \cdot p(x_n \mid y = 0)$$

- $p(x_i \mid y = 0)$ is the probability you observe $x_i$ given that an email is not spam.[2]

A more complicated method might take dependencies into account.

---

[2] Recall, $x_i$ is either 0 when word $i$ is not present, or 1 when word $i$ is present.

## Final Naive Bayes Classifier

**Training/Modeling:** Use existing data to compute:

- $p(y = 0), p(y = 1)$
- For all $i$:
  - Compute $p(0$ at position $i \mid y = 0), p(1$ at position $i \mid y_0)$
  - Compute $p(0$ at position $i \mid y = 1), p(1$ at position $i \mid y = 1)$

**Prediction:**

- For all $i$:
  - Compute $p(\mathbf{x} \mid y = 0) = \prod_i p(x_i \mid y = 0)$
  - Compute $p(\mathbf{x} \mid y = 1) = \prod_i p(x_i \mid y = 1)$
- Return

$$\arg\max \left[ p\left(\mathbf{x} \mid y = 0\right) \cdot p\left(y = 0\right), p\left(\mathbf{x} \mid y = 1\right) \cdot p\left(y = 1\right) \right].$$