# CPSC 440/540: Machine Learning

Bayesian Learning

Winter 2023

# Last Time: Monte Carlo Methods

- Monte Carlo approximates expectation of random functions:

$$\mathbb{E}[g(x)] = \sum_{x \in \mathcal{X}} g(x) p(x)$$

$$\mathbb{E}[g(x)] = \int_{x \in \mathcal{X}} g(x) p(x) \, dx$$

pmf of discrete variable X

pdf of continuous X

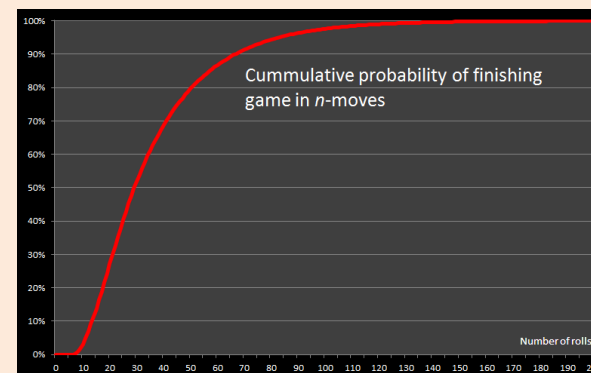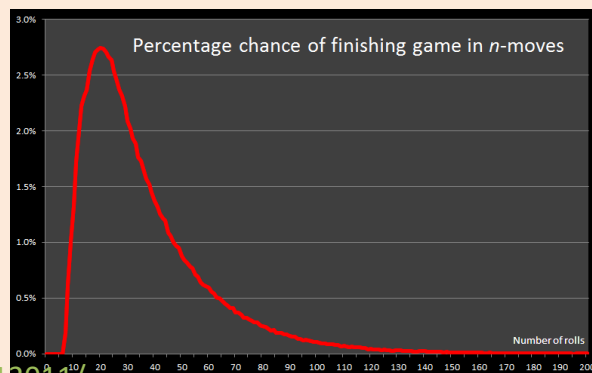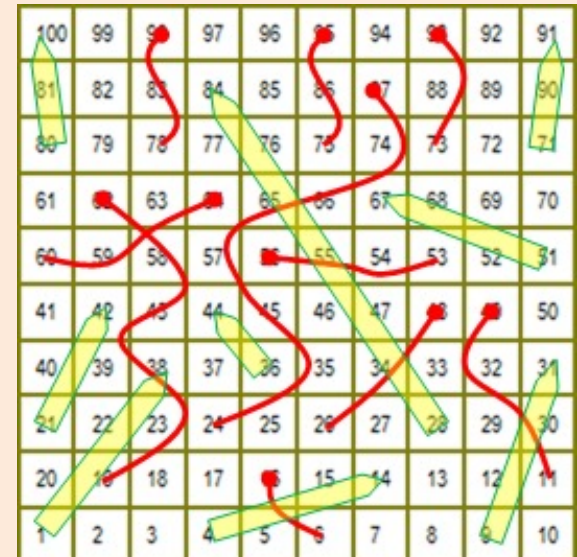- Approximation is average of function $g$ applied to samples from $p$:

$$\mathbb{E}[g(X)] \approx \frac{1}{n} \sum_{i=1}^{n} g(x^i)$$

- Can approximate a wide variety of quantities by changing $g$:
  - Mean: $g(x) = x$.
  - Probability of event 'A': $g(x) = \mathbb{1}[\text{"A happened"}]$.
  - CDF: $g(x) = \mathbb{1}[x \leq c]$.
- This is useful when:
  - You know how to sample from $p(x)$.
  - You do not know how to efficiently compute $\mathbb{E}[g(x)]$.
  - Are patient and/or don't care about being precise, because it converges slowly.

# Monte Carlo for Snakes and Ladders

bonus!

- Consider the children's game "Snakes and Ladders":
  - Start on '1', roll die, move marker, go up/down on ladder/snake, end at 100.
  - No decisions, so you can simulate the game.
- How many turns does it take for this game to end?
  - Simulate game many times, count number of turns.
  - Compute average number of turns.
- Probability and cumulative probability:

# Conditional Probabilities with Monte Carlo

- We often want to compute conditional probabilities.
  - "What is the probability that the game will go more than 100 turns, if it already went 50 turns?"

- A Monte Carlo approach for estimating $p(A \mid B)$:
  - Generate a large number of samples.
  - Use Monte Carlo estimate of $p(A, B)$ and $p(B)$ to approximate conditional:

$$p(A \mid B) = \frac{p(A, B)}{p(B)} \approx \frac{\sum_{i=1}^{\hat{n}} \mathbb{1}[\text{"A and B happened"}]}{\sum_{i=1}^{\hat{n}} \mathbb{1}[\text{"B happened"}]}$$

  - Frequency of the first event, in samples consistent with the second event.
    - This is the MLE for a binary variable that is 1 when A happens, conditioned on B happening.

- This is a special case of rejection sampling (general case later).
  - Unfortunately, if B is rare then most samples are "rejected" (ignored).
  - The conditional probability demo here has a good visualization of this.

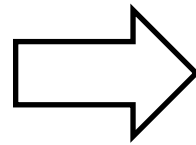# Next Topic: MLE and MAP for Categorical

# MLE for Categorical Distribution

- Now we will consider how to train a categorical model ("learning").
  - Goal is to go from samples to an estimate of parameters $\theta_1, \theta_2, \dots, \theta_k$:

| Party? |
|--------|
| LIB |
| CPC |
| NDP |
| LIB |
| GRN |

X =

⟹ p(x = LIB) = 0.34, p(x=NDP) = 0.34,
p(x = CPC) = 0.27, p(x=GRN) = 0.03,
p(x = PPC) = 0.02.

- As before we will first consider maximum likelihood estimation:

$$\hat{\widehat{\Theta}} \in argmax\{ p(x^1, x^2, \dots, x^n \mid \widehat{\Theta}) \} \rightarrow \{\theta_1, \theta_2, \dots, \theta_k\}$$

  - In this case the MLE is given by $\theta_c = \frac{n_c}{n}$ ($n_c$ is number 'c' examples).
    - If "34% of your samples are LIB, your guess for $\theta_{LIB}$=0.34".
    - As with Bernoulli, the derivation of the MLE is not as a simple as the result.

# Derivation of MLE (that does not work)

*bonus!*

- Last time we showed that the likelihood has the form:

$$p(X \mid \Theta) = \Theta_1^{n_1} \Theta_2^{n_2} \cdots \Theta_k^{n_k}$$

$$x^1, x^2, \ldots, x^n \qquad \Theta_1, \Theta_2, \ldots, \Theta_k$$

- Let's take the log:

$$\log p(X \mid \Theta) = n_1 \log \Theta_1 + n_2 \log \Theta_2 + \cdots n_k \log \Theta_k$$

- Take the derivative for a particular $\theta_c$:

$$\nabla_{\theta_c} \log p(X \mid \Theta) = \frac{n_c}{\theta_c}$$

- Set derivative equal to zero:

$$0 = \frac{n_c}{\theta_c}$$

- ...huh?

# Derivation of MLE: Handling "Sum to 1"

- "Set derivative of log-likelihood equal to 0" does not work.
  - Because of constraint that the $\theta_c$ must sum to 1, derivative is not zero at MLE.

- Approaches used in textbooks to enforce constraints:
  - Use "Lagrange multipliers" and find stationary point of "Lagrangian".
  - Define $\theta_k = 1 - \sum_{c=1}^{k-1} \theta_c$ to make it unconstrained.
  - See StackExchange thread [here](#).

- We will take a different approach to make it unconstrained:
  1. Use a unnormalized parameterization $\tilde{\theta}_c$ that doesn't have constraints.
  2. Compute the MLE for the $\tilde{\theta}_c$ by setting log-likelihood derivative to zero.
  3. Convert from the $\tilde{\theta}_c$ parameters to our usual $\theta_c$ parameters by normalizing.

# Unconstrained Parameterization

- Consider categorical distribution with unnormalized parameters:

$$p\left(x = c \mid \tilde{\theta}_1, \tilde{\theta}_2, \ldots, \tilde{\theta}_k\right) \propto \tilde{\theta}_c$$

  - To give non-negative probabilities, we require that $\bar{\theta}_c \geq 0$ for all 'c'.
- The normalized probability can then be written:

$$p\left(x = c \mid \tilde{\theta}\right) = \frac{\tilde{\theta}_c}{\sum_{c'=1}^{k} \tilde{\theta}_c} = \frac{\tilde{\theta}_c}{Z}$$

$$Z = \sum_{c'=1}^{k} \tilde{\theta}_c \text{ is called the "normalizing constant"}$$

  - The "normalizing constant" makes the probability sum to 1 across $c$ values.
    - So we do not need to an explicit "sum to 1" constraint.
  - We convert from unnormalized to normalized by dividing by $Z$: $\theta_c = \frac{\tilde{\theta}_c}{Z}$.

It is constant in terms of 'x' but a function of $\tilde{\theta}_1, \tilde{\theta}_2, \ldots, \tilde{\theta}_k$

# Derivation of MLE (that does work)

- Using the unnormalized parameters in the likelihood gives:

$$p(X \mid \Theta) = \left(\frac{\tilde{\theta_1}}{Z}\right)^{n_1} \left(\frac{\tilde{\theta_2}}{Z}\right)^{n_k} \cdots \left(\frac{\theta_K}{Z}\right)^{n_K} = \frac{\tilde{\theta_1}^{n_1} \tilde{\theta_2}^{n_k} \cdots \theta_K^{n_K}}{Z^n}$$

- Let's take the log: $\log p(X \mid \Theta) = n_1 \log(\hat{\theta_1}) + n_2 \log(\hat{\theta_2}) + \cdots + n_k \log(\tilde{\theta_K}) - n \log Z$

- Take the derivative for a particular $\theta_c$: $\nabla_{\hat{\theta_c}} p(X \mid \Theta) = \frac{n_c}{\tilde{\theta_c}} - \frac{n}{Z}$

- Set derivative equal to zero: $0 = \frac{n_c}{\tilde{\theta_c}} - \frac{n}{Z}$

- Solve for $\tilde{\theta_c}$: $\frac{\tilde{\theta_c}}{Z} = \frac{n_c}{n}$ $\Rightarrow$ Convert to normalized: $\theta_c = \frac{n_c}{n}$

(and possible to show this maximizes likelihood)

$$\log Z = \log\left(\sum_{c=1}^{K} \tilde{\theta_c}\right)$$

$$\nabla_{\tilde{\theta_c}} \log Z = \frac{1}{\sum_{c=1}^{K} \tilde{\theta_c}} = \frac{1}{Z}$$

# MAP Estimation and Dirichlet Prior

- As before, we may prefer to use a MAP estimate over the MLE.
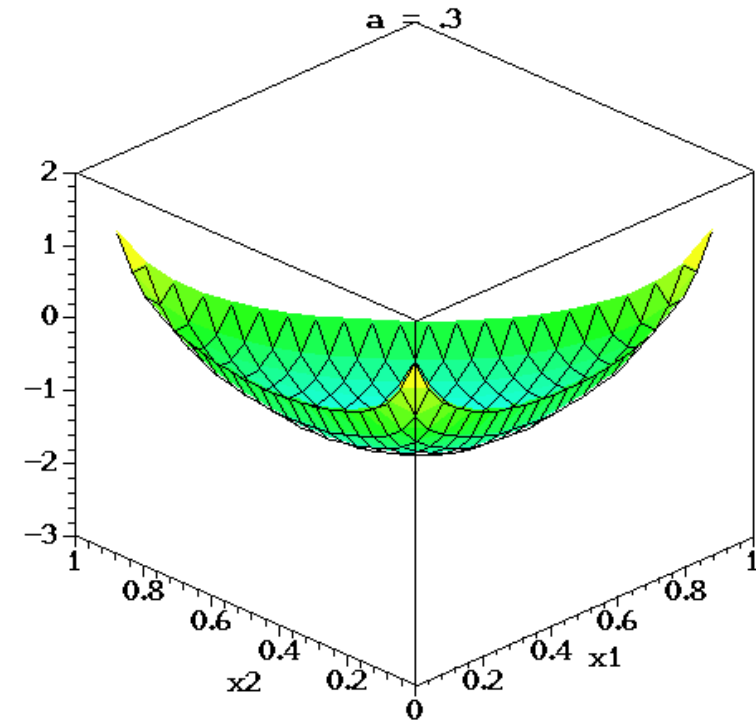  - Often becomes more important as $k$ grows.
    - More parameters to [over]fit.

- Most common prior for categorical is the Dirichlet distribution:

$$p(\Theta_1, \Theta_2, \dots, \Theta_k \mid \alpha_1, \alpha_2, \dots, \alpha_k) \propto \Theta_1^{\alpha_1 - 1} \Theta_2^{\alpha_2 - 1} \cdots \Theta_k^{\alpha_k - 1}$$

  - Generalization of the beta distribution to $k$ classes (requires $\alpha_c > 0$).

- This is a distribution over $\Theta$ values:
  - Since the $\Theta$ parameterize probabilities,
    Dirichlet is a probability distribution over possible probability distributions.

# Dirichlet Distribution

- Wikipedia's visualizations of Dirichlet distribution for k=3:



- Can bias towards various types of probabilities.

# MAP Estimation and Dirichlet Prior

- The MAP for a categorical with Dirichlet prior is given by:

$$\hat{\theta}_c = \frac{n_c + \alpha_c - 1}{\sum_{c'=1}^{k}\left[n_{c'} + \alpha_{c'} - 1\right]}$$

  – Derivation is similar to the MLE derivation.

- Dirichlet has $k$ hyper-parameters $\alpha_c$.

  – We often set $\alpha_c = \alpha$ for some constant $\alpha$ (reduces to 1 hyper-parameter).

  – This simplifies the MLE to:

$$\hat{\theta}_c = \frac{n_c + \alpha - 1}{\sum_{c'=1}^{k} n_{c'} + K(\alpha - 1)}$$

  – With $\alpha = 2$, we get Laplace smoothing ("add 1 to count of each class").

# Posterior for Categorical Likelihood + Dirichlet Prior

- People use the Dirichlet because posterior has a simple form:

$$p(\Theta \mid X, \alpha) \propto p(X \mid \Theta) p(\Theta \mid \alpha) \propto \theta_1^{n_1} \theta_2^{n_2} \cdots \theta_k^{n_k} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \cdots \theta_k^{\alpha_k-1}$$

$$\{\theta_1, \theta_2, \cdots, \theta_k\} \qquad (\alpha_1, \alpha_2, \cdots, \alpha_n)$$

Assuming data is independent of parameters given hyper-parameters

$$= \theta_1^{(n_1+\alpha_1)-1} \theta_2^{(n_2+\alpha_2)-1} \cdots \theta_k^{(n_k+\alpha_k)-1}$$

$$= \theta_1^{\tilde{\alpha}_1-1} \theta_2^{\tilde{\alpha}_2-1} \cdots \theta_k^{\tilde{\alpha}_k-1}$$

- This is another Dirichlet distribution with "updated" parameters $\tilde{\alpha}_c$.

  - Where $\tilde{\alpha}_c = n_c + \alpha_c$.

  - Again, make sure you understand why we can recognize this as a Dirichlet.

    - The normalizing constant must be the normalizing constant for the Dirichlet.

$$Z = \int_0^1 \int_0^1 \cdots \int_0^1 \theta_1^{\tilde{\alpha}_1-1} \theta_2^{\tilde{\alpha}_2-1} \cdots \theta_k^{\tilde{\alpha}_k-1} \, d\theta_1 \, d\theta_2 \cdots d\theta_k$$

# Conjugate Priors

- We have now some examples of a convenient phenomenon:
  - If we put a beta prior on a Bernoulli likelihood, posterior is beta.
    - Same happens if you put beta prior on binomial/geometric: posterior is beta.
  - If we put a Dirichlet prior on a categorical likelihood, posterior is Dirichlet.

- In these situations, we say the prior is conjugate to the likelihood.
  - With conjugate priors, the prior and posterior come from the same "family".

$$x \sim D(\theta), \quad \theta \sim P(\lambda) \quad \Rightarrow \quad \theta \mid x \sim P(\lambda')$$

this means "has the probability distribution of"

  - The posterior will look like the prior with "updated" parameters.

- Many computations become easier when we use conjugate priors.
  - Because we have an explicit formula for the posterior distribution.
  - But not all distributions have conjugate priors.

# Next Topic: Bayesian Learning

# Problems with MAP

- With good hyper-parameters, MAP usually outperforms MLE.

- But MAP is still weird.
    - Recall that we said that decoding the mode can do weird things.
        - The value with highest probability/PDF may not represent "typical" behavior.
    - MAP is *maximum a posteriori*, the posterior mode.

- MAP is fine if you want to find parameters with highest probability, but in ML usually the goal is to make predictions (or decisions).
    - Our ultimate goal is not just to find the best parameters.

- You can show that MAP is a sub-optimal way to make predictions.

# Example: "Two Heads" with "Fair vs. Unfair" Prior

- Suppose you have a Bernoulli variable and the following prior:
  - $p(\theta = 0.5) = 0.5$ and $p(\theta = 1) = 0.5$.
    - You think coin has 50% chance of being fair, 50% chance of "always landing head".
- The first two coin flips are "head".
  - $x^1 = 1$, $x^2 = 1$.
- What is the probability that the third flip will be a "head"?
  - MAP approach:

  $$1.\ \text{Find}\ \hat{\theta} \in \arg\max\{p(\theta \mid X)\} \equiv \arg\max\{p(X \mid \theta)p(\theta)\}$$

  $$2.\ \text{Compute}\ p(x^3 = 1 \mid \hat{\theta} = 1) = 1$$

  $$\theta = \tfrac{1}{2}$$
  $$\theta = 1$$
  $$(\tfrac{1}{2})(\tfrac{1}{2})(\tfrac{1}{2}) = \tfrac{1}{8} \qquad (1)(1)(\tfrac{1}{2}) = \tfrac{1}{2}$$

  - MAP predicts 100% chance of head.
    - But the MAP "decoding" of the parameters is over-confident.
      - There was a 1/4 chance of seeing two heads from the fair coin.

  $$\text{Since}\ \tfrac{1}{2} > \tfrac{1}{8},\ \text{set}\ \hat{\theta} = 1$$

# Example: "Two Heads" with "Fair vs. Unfair" Prior

- Can compute correct probability using marginalization rule over $\theta$:

$$p(x^3=1 \mid X) = \sum_{\theta \in \{0.5, 1\}} p(x^3=1, \theta \mid X) = \sum_{\theta \in \{0.5, 1\}} p(x^3=1 \mid \theta, X) \, p(\theta \mid X)$$

the probability we want

marg. rule

product rule

prediction given $\theta$

posterior

- The correct probability weights possible predictions by posterior.
  - Assume $x^3$ is independent of $X$ once we know $\theta$: $\quad p(x^3=1 \mid \theta, X) = p(x^3=1 \mid \theta)$
  - Use Bayes rule to compute posterior and get final answer:

$$p(\theta \mid X) = \frac{p(X \mid \theta) p(\theta)}{\sum_{\theta'} p(X \mid \theta') p(\theta')}$$

$\theta = \frac{1}{2} \rightarrow \quad \dfrac{1/8}{1/8 + 1/2} = \dfrac{1}{5}$

$\theta = 1 \rightarrow \quad \dfrac{1/2}{1/8 + 1/2} = \dfrac{4}{5}$

plug in

Probability from "fair" case    Probability from "unfair" case

$$p(x^3=1 \mid X) = \left(\tfrac{1}{2}\right) \cdot \left(\tfrac{1}{5}\right) + (1) \left(\tfrac{4}{5}\right)$$

$$= \frac{9}{10}$$

# Bayesian Approach to Machine Learning

- MAP predicted 100% chance that third coin would be a head.
  - But the correct value was only 90% (obtained by marginalizing over $\theta$).

- "Compute correct probability by marginalizing over parameters" is called the Bayesian approach to machine learning.
  - MAP approach optimizes posterior over parameter values.
    - Searches for the single "best" parameter value according to posterior.
  - Bayesian approach marginalizes posterior over parameter values.
    - Considers all possible parameter values, but upweighting ones with high posterior.

- MAP and Bayes are similar if posterior is "concentrated" at one $\theta$.
  - But if there are many reasonable $\theta$, Bayes can be much better.

# Digression: Review of Independence

- Let $A$ and $B$ be random variables taking values $a \in \mathcal{A}$ and $b \in \mathcal{B}$.

- We say that $A$ and $B$ are independent if for all $a$ and $b$ we have:

$$p(a,b) = p(a)p(b)$$

- To denote independence of $A$ and $B$ we often use the notation:

$$A \perp\!\!\!\perp B$$

- The product of Bernoullis model assumes mutual independence:

$$X_i \perp\!\!\!\perp X_j \quad \text{for all } i \text{ and } j$$

this is the "mutual" part

# Digression: Review of Independence

- For independent $A$ and $B$ we have:

$$p(a \mid b) = \frac{p(a, b)}{p(b)} = \frac{p(a) p(b)}{p(b)} = p(a)$$

- We can also use this as a more intuitive definition:
  - $A$ and $B$ are independent if for all $a$ and $b$ where $p(b) \neq 0$ we have:

$$p(a \mid b) = p(a)$$

    - In words: "knowing $b$ tells us nothing about $a$" (and vice versa: p(b | a)=p(b)).
    - This will often simplify calculations.

- Useful fact that can help determine if variables are independent:
  - $A \perp\!\!\!\perp B$ iff $p(a, b) = f(a) g(b)$ for some functions $f$ and $g$.

# Digression: Review of Conditional Independence

- We say that $A$ is conditionally independent of $B$ given $C$ if:

$$p(a, b \mid c) = p(a \mid c) p(b \mid c) \quad \text{for all } a, b, \text{ and } c \text{ with } p(c) \neq 0$$

  – Same as independence definition, but "knowing extra stuff" $C$.

- Or, alternatively:

$$p(a \mid b, c) = p(a \mid c) \quad \text{or} \quad p(b \mid a, c) = p(b \mid c)$$

  – "If you know $C$, then *also* knowing $B$ would tell you nothing about $A$."

- We often write this as: $A \perp\!\!\!\perp B \mid C$

- In naïve Bayes we assume $X_i \perp\!\!\!\perp X_j \mid Y$ for all i and j.

  – As we saw, this makes inference and learning easy.

# Standard ML Independence Assumptions (MEMORIZE)

- In machine learning we typically make a standard set of independence assumptions:
  - IID assumption: training examples are independent of each other.

$$x^i \perp\!\!\!\perp x^j$$

  - "If you see example $x^i$, it doesn't tell you anything about $x^j$."
  - Maybe better framing is $x^i \perp\!\!\!\perp x^k \mid \mathcal{D}$: they're conditionally independent given the hidden "data-generating process" $\mathcal{D}$.
  - Independence of data given parameters.

$$x^i \perp\!\!\!\perp x^j \mid \Theta$$

  - "If we know the parameters, the examples are independent of each other"
  - Again, maybe better to think of this as $x^i \perp\!\!\!\perp x^k \mid \theta, \mathcal{D}$.
  - Independence of features $X$ and parameters $w$ in **discriminative** models.

$$w \perp\!\!\!\perp X$$

  - Discriminative models assume parameters are fixed, and $w$ just transforms them to $y$ (knowing $X$ without $y$ tells you nothing).
  - Conditional independence of data and hyper-parameters, given parameters:

$$X \perp\!\!\!\perp \alpha, \beta \mid \Theta$$

  - "Given the parameters, the hyper-parameters don't tell you anything more about the data."
- Later we'll discuss the models that lead to these assumptions, and testing independence in a model.

# Bayesian Approach for Bernoulli-Beta Model

- Consider probability that $x^3 = 1$ after $x^1 = 1$ and $x^2 = 1$ with beta prior:

$$p(x^3 = 1 \mid X, \alpha, \beta) = \int_\Theta p(x^3 = 1, \Theta \mid X, \alpha, \beta) \, d\Theta \qquad \text{(marginalization rule)}$$

"posterior predictive"

$$= \int_\Theta p(x^3 = 1 \mid \Theta, X, \alpha, \beta) \, p(\Theta \mid X, \alpha, \beta) \, d\Theta \qquad \text{(product rule)}$$

$$= \int_\Theta p(x^3 = 1 \mid \Theta) \, p(\Theta \mid X, \alpha, \beta) \, d\Theta \qquad \text{(conditional independence)}$$

"prediction"    "posterior"

- Now use that posterior is a beta with parameters $\tilde{\alpha}$ and $\tilde{\beta}$.

$$= \int_\Theta \Theta \, p_\beta(\Theta \mid \tilde{\alpha}, \tilde{\beta}) \, d\Theta \qquad \text{(definition of Bernoulli and form of posterior)}$$

$$= \mathbb{E}[\Theta] \qquad \text{(expected value of } \Theta \text{ under posterior distribution)}$$

$$= \frac{\tilde{\alpha}}{\tilde{\alpha} + \tilde{\beta}} \qquad \text{(formula for expected value of } \Theta \text{ under beta)}$$

# Bayesian Approach for Bernoulli-Beta Model

- The correct probability of seeing a "head" after 2 flips in Bernoulli-beta:

$$p(x^3 = 1 \mid X, \alpha, \beta) = \int_0^1 p(x^3 = 1, \theta \mid X, \alpha, \beta) \, d\theta$$

$$= \frac{\tilde{\alpha}}{\tilde{\alpha} + \tilde{\beta}} \quad (\text{last slide})$$

$$= \frac{n_1 + \alpha}{(n_1 + \alpha) + (n_0 + \beta)}$$

- With a uniform prior, ($\alpha = \beta = 1$), then Pr($x^3 = 1 \mid x^1=1, x^2=1, \alpha, \beta$) = ¾.
  - The MAP under a uniform prior (which is MLE) would be $\theta = 1$.
    - It is less confident than MAP since it considers all possible $\theta$ values, not just the most likely.
    - Bayesian estimate is not degenerate even under a uniform prior here.
- Looks like Laplace smoothing, but trusts data less for same $\alpha$ and $\beta$.
  - For other models, MAP and Bayes can be much more different.

# Effect of Prior in Bernoulli-Beta

- In Bayesian approach, hyper-parameters $\alpha$ and $\beta$ can be thought of as "pseudo-counts".
  - The number of 0 and 1 outcomes you have in your imagination before you see any data.

- If we see 3 "heads" ($x^1=1, x^2=1, x^3=1$), the probability of a 4th under different priors:
  - Beta(1,1) prior is like seeing 1 imaginary head and 1 tail before flipping.
    - Probability is 4/5, even though all $\theta$ values under this uniform prior "equally likely".
  - Beta(3,3) prior is like seeing 3 imaginary heads and 3 tails.
    - Probability is 0.667. This is a stronger bias towards 0.5.
  - Beta(100,1) prior is like seeing 100 imaginary heads and 1 tail.
    - Probability is 0.990. This is a strong bias towards high $\theta$ values.
  - Beta(0.01,0.01) prior biases towards having an unfair coin (head or tail).
    - Probability is 0.997.

- We might hope to use an "uninformative" prior to not bias results.
  - We saw that with the "uniform" prior, Beta(1,1), it biases towards 0.5.
  - See bonus for additional details on why "uninformative" can be hard/ambiguous/impossible/undesirable.

# Motivation: Controlling Complexity

- For many application, we need complicated models.
- But complex models can overfit.
- So what should we do?

- In CPSC 340 we saw two ways to reduce overfitting:
  - Model averaging (like in random forests).
  - Regularization (like in L2-regularized linear regression).

- Bayesian methods combine both of these.
  - Average over "models", weighted by posterior (which includes regularizer).
    - Recall that the regularizer corresponds to the negative logarithm of the prior.
  - This can allow you fit extremely complicated models without overfitting.

# MAP vs Bayes for Categorical-Dirichlet

- MAP (regularized optimization) approach maximizes over parameters:

$$\hat{\Theta} \in \underset{\Theta}{\arg\max} \left\{ p(\Theta \mid X) \right\}$$

$$\equiv \underset{\Theta}{\arg\max} \left\{ p(X \mid \Theta) p(\Theta) \right\} \quad (\text{Bayes' rule})$$

$$p(x=c \mid \hat{\Theta}) = \hat{\Theta}_c$$

*(I'm not explicitly including the conditioning on the hyper-parameters $\alpha$)*

- Bayesian approach predicts by integrating over possible parameters:

$$p(x=c \mid X) = \int_{\Theta_1} \int_{\Theta_2} \cdots \int_{\Theta_K} p(x=c, \Theta \mid X) \, d\Theta_K \, d\Theta_{K-1} \cdots d\Theta_1 \quad (\text{marg. rule})$$

$$= \int_{\Theta_1} \int_{\Theta_2} \cdots \int_{\Theta_K} p(x=c \mid \Theta, X) \, p(\Theta \mid X) \, d\Theta_K \, d\Theta_{K-1} \cdots d\Theta_1 \quad (\text{product rule})$$

$$= \int_{\Theta_1} \int_{\Theta_2} \cdots \int_{\Theta_K} \Theta_c \, p(\Theta \mid X) \, d\Theta_K \, d\Theta_{K-1} \cdots d\Theta_1 \quad \begin{array}{l}(\text{independence of data}\\ \text{given parameters})\end{array}$$

- Considers all possible $\Theta$, and weights prediction by posterior for $\Theta$.
  - Posterior contains a regularizer, so this is averaging and regularizing.

$$\rightarrow \mathbb{E}(\Theta_c) \quad \begin{array}{l}(\text{mean of Dirichlet}\\ \text{posterior})\end{array}$$

# Ingredients of Bayesian Inference (MEMORIZE)

1. **Likelihood** $p(X \mid \Theta)$
   - Probability of seeing data given parameters.

2. **Prior** $p(\Theta \mid A)$.
   - Belief that parameters are correct before we have seen data.

3. **Posterior** $p(\Theta \mid \boldsymbol{X}, A)$.
   - Probability that parameters are correct after we have seen data.
   - MAP maximizes, but Bayesian approach uses the whole distribution.

4. **Posterior predictive** $p(\tilde{X} \mid \boldsymbol{X}, A)$ **(NEW)**.
   - Probability of new data $\tilde{X}$ given old data $\boldsymbol{X}$, integrating over parameters.
     - Specifically, we average the likelihood of $\tilde{X}$, weighted by the posterior of $\theta$ given $\boldsymbol{X}$.
   - Bayesian approach uses this distribution for inference.

# Bayesian Approach: Discussion

- Our previous "learn then predict" approaches (MLE and MAP):
  - Optimize parameters $\theta$ (learning).
  - Do inference with the parameter estimate $\hat{\theta}$ (inference).

- Bayesian approach doesn't really have a separate "learning phase".
  - There is <span style="color:red">no optimization</span> of the parameter $\theta$.
  - You just skip to doing <span style="color:green">inference with the posterior predictive</span>.
    - Consider all parameters $\theta$.

- In practice, it often still looks like "learn then predict".
  - Characterize the form of the posterior ("learning").
  - Make predictions by doing integrals with the posterior (inference).

# Bayesian Approach: Discussion

- The Bayesian approach is the optimal way to use a probabilistic model.
  - It's what the rules of probability say we should do.
  - …if you believe in your probability model (prior + likelihood).

- If the prior is bad, Bayesian approach can be harmful.
  - Bayesian approach historically criticized since it requires "subjective" prior.
  - But all models are based on "subjective" assumptions, sometime hidden!

- As we see more data, Bayesian posterior concentrates on MLE.
  - MLE/MAP/Bayes usually more or less agree for large datasets.

- Real problem with the Bayesian approach is that integrals are hard.
  - Posterior and posterior predictive only have a nice form with conjugate priors.
    - Otherwise, you need to use methods like Monte Carlo or "variational" methods for inference.

# Uninformative Priors and Jeffreys Priors

*bonus!*

- We might want to use an uninformative prior to not bias results.
  - But this is often hard/impossible to do.

- We might think the uniform distribution, Beta(1,1), is uninformative.
  - But posterior will be biased towards 0.5 compared to MLE.
  - And if you use a different parameterization it won't stay uniform.

- We might think to use "pseudo-count" of 0, Beta(0,0), as uninformative.
  - But posterior isn't a probability until we see at least one head and one tail.

- Some argue that the "correct" uninformative prior is Beta(0.5,0.5).
  - This prior is invariant to the parameterization, which is called a Jeffreys prior.

# Summary

- **MLE for categorical** distribution:
  - Write using **unnormalized** parameters and **normalizing constant** 'Z'.
- **Dirichlet distribution**:
  - "Probability distribution over discrete probability distributions".
  - When used as prior for categorical, posterior is also Dirichlet.
  - MAP estimate with Dirichlet prior gives generalization of Laplace smoothing.
- **Conjugate prior**:
  - Prior for a particular likelihood such that posterior is in same "family".

- **Conditional independence** of A and B [given C].
  - "Knowing A tells you nothing about B [if you also know C]".
  - Independence assumptions often simplify computations.
  - In ML we make a **standard set of independence assumptions**.
    - Data and hyper-parameters are independent given parameters.
- **Bayesian learning**.
  - Do inference with the **posterior predictive** (no "learning" phase).
  - Can be viewed as regularizing and averaging over parameters (harder to overfit).
  - Involves solving unpleasant integrals (unless you have a conjugate prior).

- Next time: priors on priors + relaxing IID.