

Lecture 6: MAP Estimation

Priors · Posteriors · Regularization Connection

Previously, on Lecture 5...

MoM: Set population moments = sample moments. Simple but can give impossible values.

MLE: $\hat{\theta} = \arg \max \ell(\theta)$. Pick the θ that makes the data most probable.

Properties: Consistent, asymptotically normal, efficient, invariant.

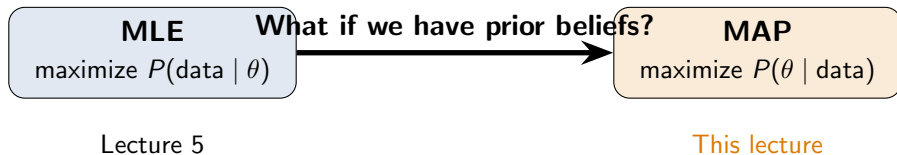
MLE = ML: Gaussian noise \rightarrow MSE loss. Bernoulli \rightarrow cross-entropy.

But: MLE can overfit with small n or flexible models.

Today: What if we have **prior knowledge** about θ ?

Can we do better than MLE by incorporating beliefs *before* seeing data?

Where We Are

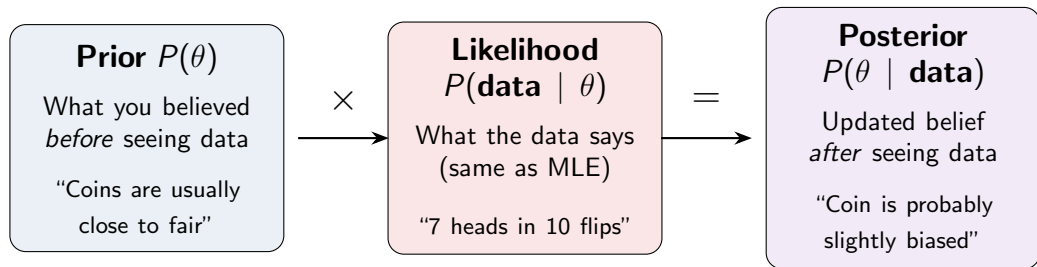


Bayes' Theorem for Parameters

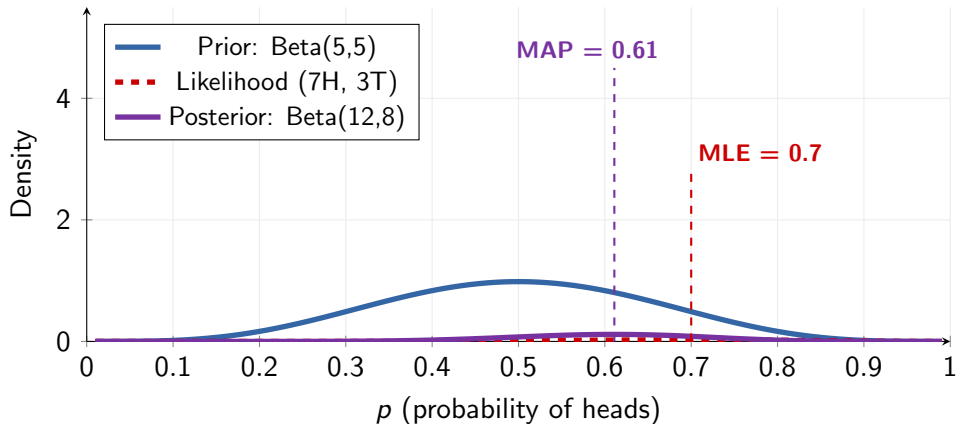
$$\underbrace{P(\theta \mid \text{data})}_{\text{posterior}} = \frac{\overbrace{P(\text{data} \mid \theta)}^{\text{likelihood}} \cdot \overbrace{P(\theta)}^{\text{prior}}}{\underbrace{P(\text{data})}_{\text{evidence}}}$$

Or simply: posterior \propto likelihood \times prior

The Three Ingredients



Visualizing the Update: Coin Bias



Prior pulls the estimate from 0.7 toward 0.5. The posterior is a **compromise**.

Conjugate Priors as Pseudo-Observations

Why did the math work out so neatly? Because Beta is **conjugate** to Binomial:

Prior: $\text{Beta}(\alpha, \beta)$

+

Data: k H, $n-k$ T

=

Posterior: $\text{Beta}(\alpha+k, \beta+n-k)$

Conjugate Priors as Pseudo-Observations

Why did the math work out so neatly? Because Beta is **conjugate** to Binomial:

$$\text{Prior: Beta}(\alpha, \beta) + \text{Data: } k \text{ H, } n-k \text{ T} = \text{Posterior: Beta}(\alpha+k, \beta+n-k)$$

The prior acts like “fake data” you’ve already seen:

$\text{Beta}(\alpha, \beta)$ = pretend you already observed $\alpha-1$ heads and $\beta-1$ tails.

$\text{Beta}(5, 5)$: “I’ve seen 4H and 4T” (8 pseudo-observations).

After 7H, 3T (10 real obs): posterior = $\text{Beta}(12, 8)$ = “11H, 7T out of 18 total.”

As $n \rightarrow \infty$, the pseudo-observations become negligible \Rightarrow posterior \rightarrow likelihood \Rightarrow MAP \rightarrow MLE.

MAP = Mode of the Posterior

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P(\theta \mid \text{data}) = \arg \max_{\theta} [\ell(\theta) + \log P(\theta)]$$

Maximize: log-likelihood + log-prior

MLE: $\arg \max_{\theta} \ell(\theta)$

+ prior

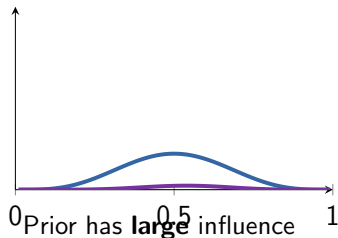


MAP: $\arg \max_{\theta} \ell(\theta) + \log P(\theta)$

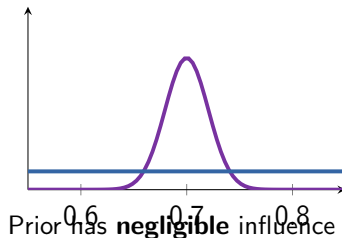
MAP = MLE with an extra penalty/bonus term from the prior.

When Does the Prior Matter?

Small n (e.g., $n = 5$)



Large n (e.g., $n = 500$)



With enough data, the likelihood dominates \Rightarrow MAP \approx MLE.
The prior is “washed out” by the data.

The Key Connection: Regularization = MAP

$$\text{MAP: } \hat{\theta} = \arg \max_{\theta} [\ell(\theta) + \log P(\theta)]$$

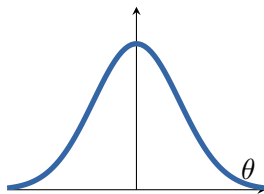
is the same as

$$\text{Regularization: } \hat{\theta} = \arg \min_{\theta} [-\ell(\theta) + \lambda \cdot \text{penalty}(\theta)]$$

The log-prior acts as a **penalty on the parameters**.

Gaussian Prior \Leftrightarrow Ridge (L2) Regression

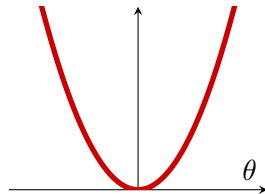
Gaussian prior



$$P(\theta) = \mathcal{N}(0, \tau^2)$$



L2 penalty



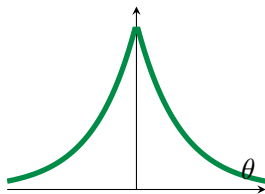
$$-\log P(\theta) \propto \|\theta\|_2^2$$

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \left[\sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \theta)^2 + \lambda \|\theta\|_2^2 \right]$$

This is exactly **Ridge regression**! $\lambda = \sigma^2 / \tau^2$

Laplace Prior \Leftrightarrow Lasso (L1) Regression

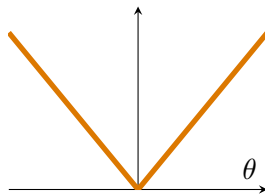
Laplace prior



$$P(\theta) \propto e^{-|\theta|/b}$$



L1 penalty



$$-\log P(\theta) \propto \|\theta\|_1$$

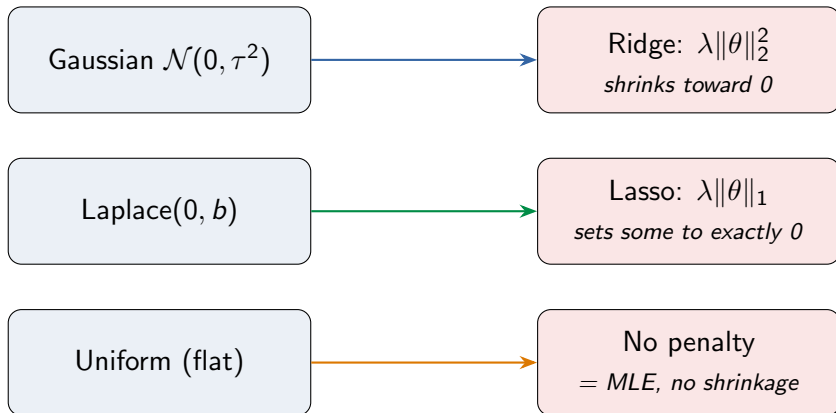
$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} [\sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \theta)^2 + \lambda \|\theta\|_1]$$

This is exactly **Lasso regression**! Encourages **sparse** solutions ($\theta_j = 0$).

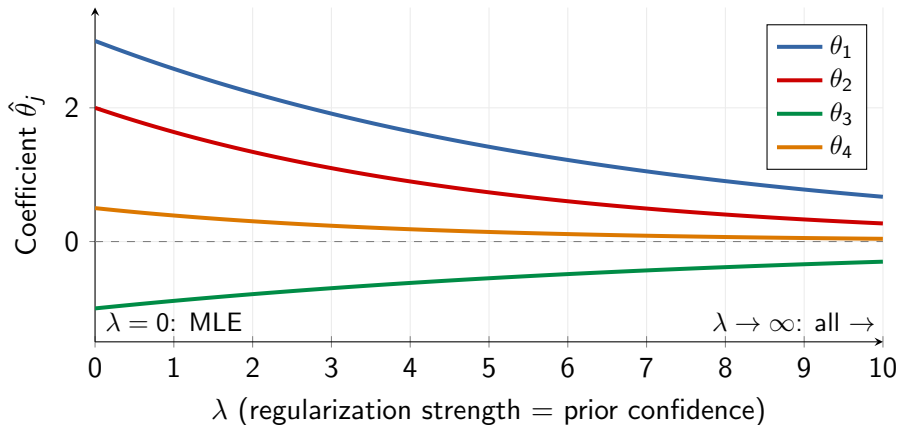
The Regularization Map

Prior (Bayesian)

Penalty (Frequentist)



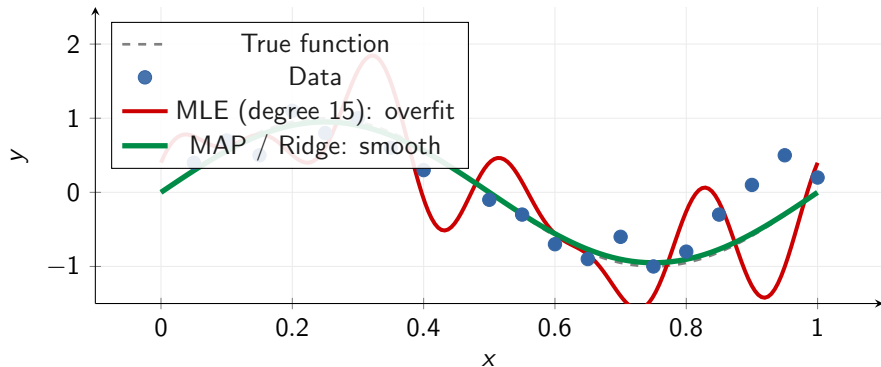
Visualizing Ridge Shrinkage



Increasing λ = stronger prior = more shrinkage = less overfitting (but more bias).

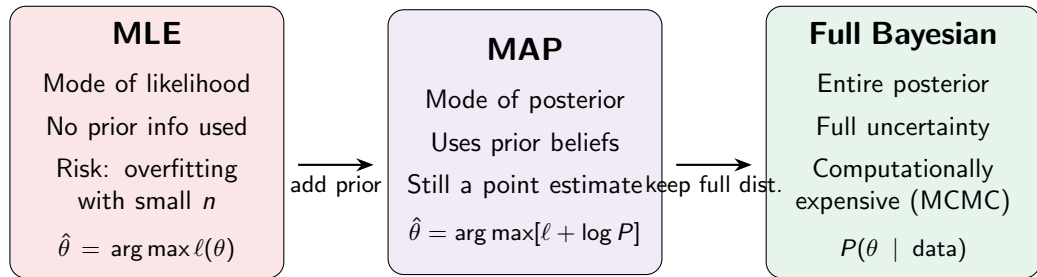
MLE vs MAP: The Overfitting Story

Fit a polynomial to noisy data. MLE uses all parameters freely; MAP penalizes large coefficients.



The prior says “coefficients should be small” \Rightarrow smoother fit \Rightarrow better generalization.

Three Philosophies



When to Use What

MLE when:

- Large n (prior doesn't matter)
- No reliable prior info
- Simplicity is valued

MAP when:

- Small n (need regularization)
- Have domain knowledge
- Want a point estimate fast

Full Bayesian when:

- Uncertainty quantification matters (medical, safety)
- Model comparison needed
- Computational cost is acceptable

Practical: Priors and Posteriors

1. **Coin bias estimation:**

- ▶ Start with Beta(1,1), Beta(5,5), Beta(50,50) priors
- ▶ Observe 7 heads in 10 flips
- ▶ Plot prior, likelihood, and posterior for each
- ▶ Compare the MAP estimates — how much does the prior pull?

2. **Ridge regression as MAP:**

- ▶ Fit linear regression with $\lambda = 0, 0.1, 1, 10, 100$
- ▶ Plot coefficients vs λ (shrinkage path)
- ▶ Observe: larger λ = stronger prior = more shrinkage

3. **Visualize:** Plot the prior/likelihood/posterior for a simple 1D Normal with known σ^2 , varying the prior variance τ^2

Homework

1. For $X_1, \dots, X_n \sim N(\mu, \sigma_0^2)$ (known σ_0^2) with prior $\mu \sim N(m, \tau^2)$:
derive the MAP estimator $\hat{\mu}_{\text{MAP}}$. Show it is a weighted average of \bar{X} and m .
What happens as $\tau^2 \rightarrow \infty$? As $n \rightarrow \infty$?
2. A coin is flipped 20 times with 14 heads. Compute the MAP estimate of p under:
(a) Beta(1, 1), (b) Beta(5, 5), (c) Beta(50, 50) priors.
Compare with the MLE. Which prior has the most influence?
3. Show that Ridge regression $\hat{\theta} = \arg \min [\|y - X\theta\|^2 + \lambda \|\theta\|^2]$
has the closed-form solution $\hat{\theta} = (X^\top X + \lambda I)^{-1} X^\top y$.
Why does this always have a unique solution, even when $X^\top X$ is singular?

Questions?

Next: Sampling distributions and confidence intervals