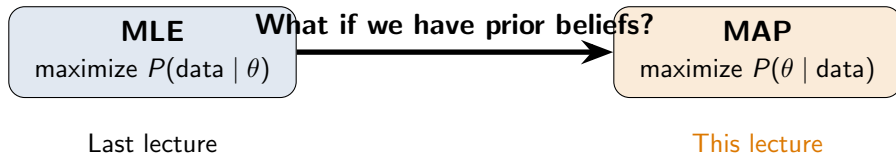


Lecture 2b: MAP Estimation

Priors, Posteriors, and the Regularization Connection

Where We Are

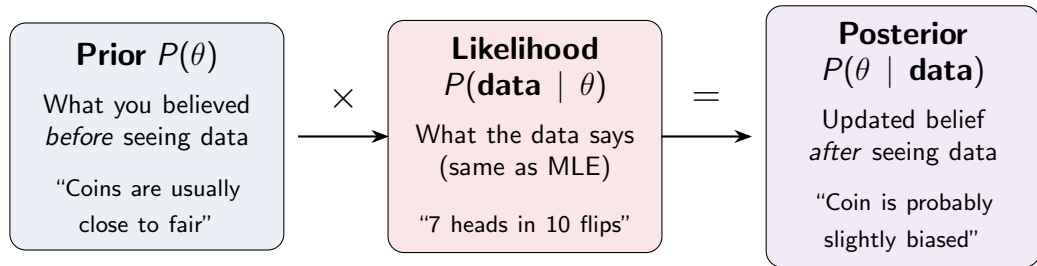


Bayes' Theorem for Parameters

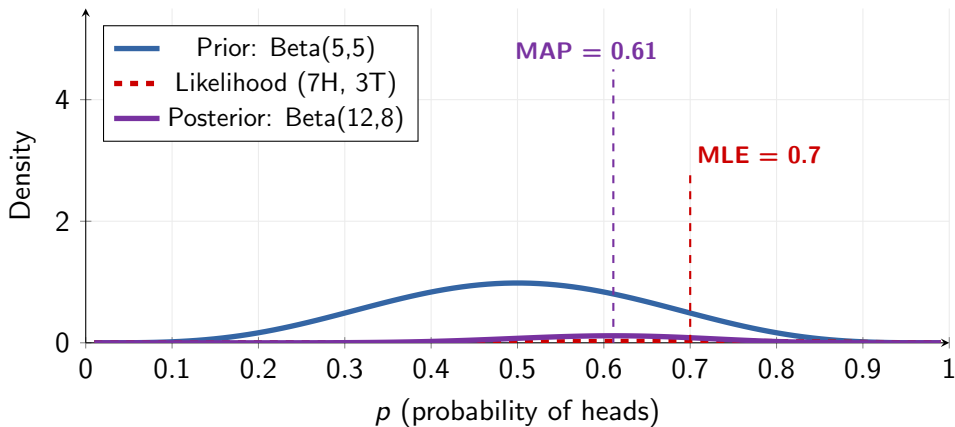
$$\underbrace{P(\theta \mid \text{data})}_{\text{posterior}} = \frac{\overbrace{P(\text{data} \mid \theta)}^{\text{likelihood}} \cdot \overbrace{P(\theta)}^{\text{prior}}}{\underbrace{P(\text{data})}_{\text{evidence}}}$$

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**.

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

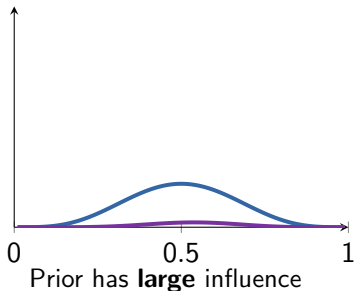


MAP: $\arg \max_{\theta} \ell(\theta) + \log P(\theta)$

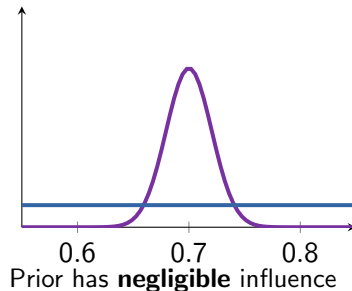
MAP = MLE with an extra penalty/bonus term from the prior.

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

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

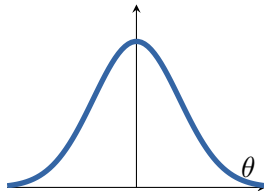
is the same as

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

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

Gaussian Prior \Leftrightarrow Ridge (L2) Regression

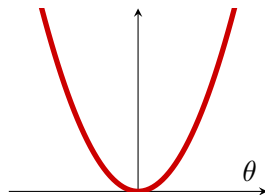
Gaussian prior



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



L2 penalty



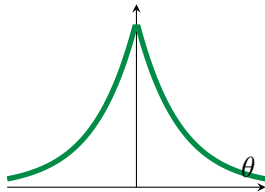
$$-\log P(\theta) \propto \|\theta\|_2^2$$

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

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

Laplace Prior \Leftrightarrow Lasso (L1) Regression

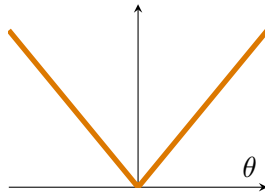
Laplace prior



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



L1 penalty



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

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

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

The Regularization Map

Prior (Bayesian)

Penalty (Frequentist)

Gaussian $\mathcal{N}(0, \tau^2)$



Ridge: $\lambda \|\theta\|_2^2$

Shrinks all

Laplace(0, b)



Lasso: $\lambda \|\theta\|_1$

Sets some

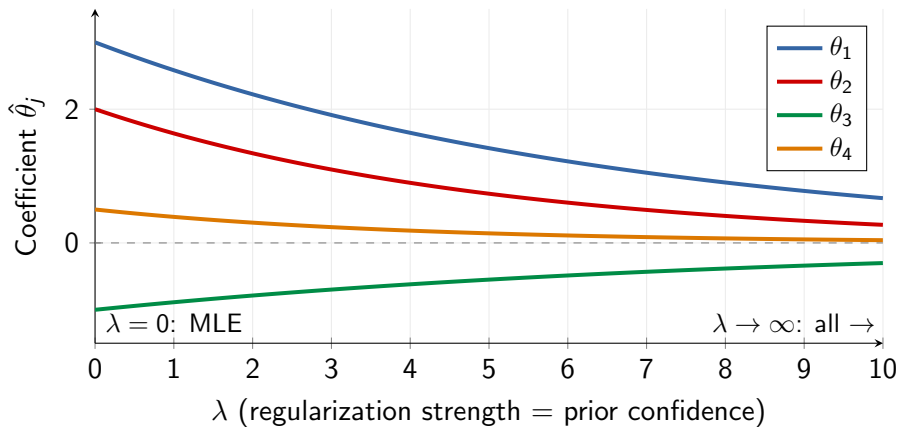
Uniform (flat)



No penalty (MLE)

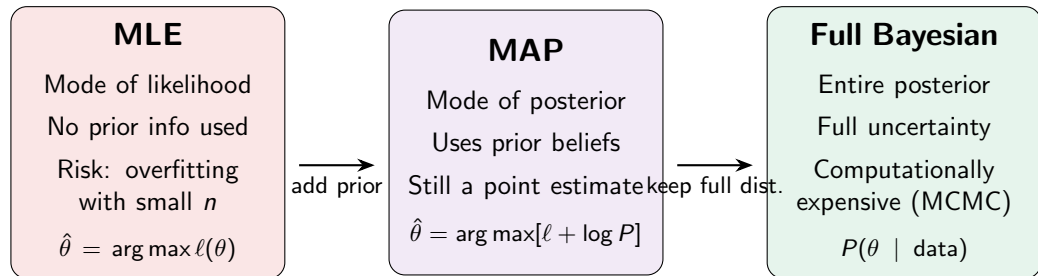
No shrinkage

Visualizing Ridge Shrinkage



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

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 is critical (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

Questions?

Next lecture: Estimator Quality — Bias, Variance, and the Tradeoff