

Байесовские Методы. Лекция 3

БИВТ-23-9/10-ИСАД

23 сентября 2025 г.

Модель

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

- X — матрица признаков размера $n \times d$.

Модель

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

- X — матрица признаков размера $n \times d$.
- θ — вектор коэффициентов (неизвестные параметры).

Модель

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

- X — матрица признаков размера $n \times d$.
- θ — вектор коэффициентов (неизвестные параметры).
- y — вектор наблюдений.

Классическая линейная регрессия

Модель

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

- X — матрица признаков размера $n \times d$.
- θ — вектор коэффициентов (неизвестные параметры).
- y — вектор наблюдений.
- Ошибки ε независимы, одинаково распределены, дисперсия σ^2 .

Цель

Найти такие параметры θ , которые лучше всего объясняют данные.

Оценка параметров (MLE)

- Правдоподобие:

$$P(y | X, \theta) = \prod_{i=1}^n \mathcal{N}(y_i | \mathbf{x}_i^\top \theta, \sigma^2).$$

Оценка параметров (MLE)

- Правдоподобие:

$$P(y | X, \theta) = \prod_{i=1}^n \mathcal{N}(y_i | x_i^\top \theta, \sigma^2).$$

- Лог-правдоподобие:

$$\ell(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|y - X\theta\|^2.$$

Оценка параметров (MLE)

- Правдоподобие:

$$P(y | X, \theta) = \prod_{i=1}^n \mathcal{N}(y_i | x_i^\top \theta, \sigma^2).$$

- Лог-правдоподобие:

$$\ell(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|y - X\theta\|^2.$$

- Максимизация эквивалентна задаче МНК:

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} \|y - X\theta\|^2.$$

Оптимизационная задача (MLE)

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} \|y - X\theta\|^2$$

- Запишем функцию ошибки (сумму квадратов):

$$J(\theta) = (y - X\theta)^\top (y - X\theta).$$

Оптимизационная задача (MLE)

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} \|y - X\theta\|^2$$

- Запишем функцию ошибки (сумму квадратов):

$$J(\theta) = (y - X\theta)^\top (y - X\theta).$$

- Берём градиент по θ :

$$\nabla_{\theta} J(\theta) = -2X^\top (y - X\theta).$$

Оптимизационная задача (MLE)

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} \|y - X\theta\|^2$$

- Запишем функцию ошибки (сумму квадратов):

$$J(\theta) = (y - X\theta)^\top (y - X\theta).$$

- Берём градиент по θ :

$$\nabla_{\theta} J(\theta) = -2X^\top (y - X\theta).$$

- Приравниваем к нулю:

$$X^\top X \hat{\theta} = X^\top y.$$

Оптимизационная задача (MLE)

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} \|y - X\theta\|^2$$

- Запишем функцию ошибки (сумму квадратов):

$$J(\theta) = (y - X\theta)^\top (y - X\theta).$$

- Берём градиент по θ :

$$\nabla_{\theta} J(\theta) = -2X^\top (y - X\theta).$$

- Приравниваем к нулю:

$$X^\top X \hat{\theta} = X^\top y.$$

- Если $X^\top X$ обратима:

$$\hat{\theta} = (X^\top X)^{-1} X^\top y.$$

Нормальные уравнения

$$X^T X \hat{\theta} = X^T y$$

- Если $X^T X$ обратима:

$$\hat{\theta} = (X^T X)^{-1} X^T y.$$

Нормальные уравнения

$$X^T X \hat{\theta} = X^T y$$

- Если $X^T X$ обратима:

$$\hat{\theta} = (X^T X)^{-1} X^T y.$$

- Если $X^T X$ вырожденная (нет обратной), решение выражается через **псевдообратную Мура-Пенроуза**:

$$\hat{\theta} = X^+ y,$$

где

$$X^+ = \begin{cases} (X^T X)^{-1} X^T, & n \geq d, \text{ ранг}(X) = d, \\ X^T (X X^T)^{-1}, & n < d, \text{ ранг}(X) = n. \end{cases}$$

Нормальные уравнения

$$X^T X \hat{\theta} = X^T y$$

- Если $X^T X$ обратима:

$$\hat{\theta} = (X^T X)^{-1} X^T y.$$

- Если $X^T X$ вырожденная (нет обратной), решение выражается через **псевдообратную Мура-Пенроуза**:

$$\hat{\theta} = X^+ y,$$

где

$$X^+ = \begin{cases} (X^T X)^{-1} X^T, & n \geq d, \text{ ранг}(X) = d, \\ X^T (X X^T)^{-1}, & n < d, \text{ ранг}(X) = n. \end{cases}$$

- Это решение минимизирует $\|y - X\theta\|^2$ и среди всех возможных выбирает θ с минимальной нормой.

Нормальные уравнения и псевдообратная

Нормальные уравнения

$$X^T X \hat{\theta} = X^T y$$

- Если $X^T X$ обратима:

$$\hat{\theta} = (X^T X)^{-1} X^T y.$$

- Если $X^T X$ вырожденная (нет обратной), решение выражается через **псевдообратную Мура-Пенроуза**:

$$\hat{\theta} = X^+ y,$$

где

$$X^+ = \begin{cases} (X^T X)^{-1} X^T, & n \geq d, \text{ ранг}(X) = d, \\ X^T (X X^T)^{-1}, & n < d, \text{ ранг}(X) = n. \end{cases}$$

- Это решение минимизирует $\|y - X\theta\|^2$ и среди всех возможных выбирает θ с минимальной нормой.

Замечание

Использование псевдообратной особенно важно при $n < d$ (мало данных, много признаков).

Ограничения классической регрессии (MLE)

- **Нет априорной информации:** модель не учитывает знаний о параметрах до наблюдений.

Ограничения классической регрессии (MLE)

- **Нет априорной информации:** модель не учитывает знаний о параметрах до наблюдений.
- **Точка вместо распределения:** $\hat{\theta}$ — одно число, без оценки неопределённости.

Ограничения классической регрессии (MLE)

- **Нет априорной информации:** модель не учитывает знаний о параметрах до наблюдений.
- **Точка вместо распределения:** $\hat{\theta}$ — одно число, без оценки неопределённости.
- **Переобучение:** при $n < d$ или при сильной корреляции признаков решение становится нестабильным.

Ограничения классической регрессии (MLE)

- **Нет априорной информации:** модель не учитывает знаний о параметрах до наблюдений.
- **Точка вместо распределения:** $\hat{\theta}$ — одно число, без оценки неопределённости.
- **Переобучение:** при $n < d$ или при сильной корреляции признаков решение становится нестабильным.
- **Нет доверительных интервалов:** предсказания не содержат информацию о «надежности».

Почему байесовская регрессия?

- Добавляем **априор** на параметры θ :

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

Почему байесовская регрессия?

- Добавляем **априор** на параметры θ :

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

- Получаем **постериорное распределение**, а не точку:

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

Почему байесовская регрессия?

- Добавляем **априор** на параметры θ :

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

- Получаем **постериорное распределение**, а не точку:

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Теперь можно делать **байесовское предсказание**:

$$P(y^* \mid x^*, D) = \int P(y^* \mid x^*, \theta) P(\theta \mid D) d\theta.$$

Байесовская линейная регрессия: модель

Наблюдения

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I).$$

Априор на параметры

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

- Правдоподобие:

$$P(y \mid X, \theta) = \mathcal{N}(y \mid X\theta, \sigma^2 I).$$

Байесовская линейная регрессия: модель

Наблюдения

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I).$$

Априор на параметры

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

- Правдоподобие:

$$P(y | X, \theta) = \mathcal{N}(y | X\theta, \sigma^2 I).$$

- Априор: гауссовский, центрирован в нуле.

Байесовская линейная регрессия: модель

Наблюдения

$$y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I).$$

Априор на параметры

$$\theta \sim \mathcal{N}(0, \tau^2 I).$$

- Правдоподобие:

$$P(y | X, \theta) = \mathcal{N}(y | X\theta, \sigma^2 I).$$

- Априор: гауссовский, центрирован в нуле.
- Постериор: по теореме Байеса

$$P(\theta | D) \propto P(y | X, \theta) P(\theta).$$

Вывод постериора для параметров

Теорема Байеса

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Правдоподобие:

$$P(y \mid X, \theta) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 \right).$$

Вывод постериора для параметров

Теорема Байеса

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Правдоподобие:

$$P(y \mid X, \theta) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 \right).$$

- Априор:

$$P(\theta) \propto \exp \left(- \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

Вывод постериора для параметров

Теорема Байеса

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Правдоподобие:

$$P(y \mid X, \theta) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 \right).$$

- Априор:

$$P(\theta) \propto \exp \left(- \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

- Постериор:

$$P(\theta \mid D) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

Вывод постериора для параметров

Теорема Байеса

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Правдоподобие:

$$P(y \mid X, \theta) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 \right).$$

- Априор:

$$P(\theta) \propto \exp \left(- \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

- Постериор:

$$P(\theta \mid D) \propto \exp \left(- \frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

- Это экспонента от квадратичной формы по $\theta \rightarrow$ многомерное нормальное распределение.

Вывод постериора для параметров

Теорема Байеса

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta).$$

- Правдоподобие:

$$P(y \mid X, \theta) \propto \exp \left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 \right).$$

- Априор:

$$P(\theta) \propto \exp \left(-\frac{1}{2\tau^2} \|\theta\|^2 \right).$$

- Постериор:

$$P(\theta \mid D) \propto \exp \left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2 \right).$$

- Это экспонента от квадратичной формы по $\theta \rightarrow$ многомерное нормальное распределение.

Результат

$$P(\theta \mid D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

Интуиция: что означают формулы?

- Σ_N — ковариация постериора.

Интуиция: что означают формулы?

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^T X$, тем меньше дисперсия → растёт уверенность.

Интуиция: что означают формулы?

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^T X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.

Интуиция: что означают формулы?

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^T X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.

Интуиция: что означают формулы?

Постериор

$$P(\theta | D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.

Интуиция: что означают формулы?

Постериор

$$P(\theta \mid D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.
 - При $\tau^2 \rightarrow \infty$ (слабый априор): $\mu_N \rightarrow \hat{\theta}_{\text{MLE}}$.

Интуиция: что означают формулы?

Постериор

$$P(\theta | D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.
 - При $\tau^2 \rightarrow \infty$ (слабый априор): $\mu_N \rightarrow \hat{\theta}_{\text{MLE}}$.
 - При малом τ^2 : μ_N тяготеет к нулю.

Интуиция: что означают формулы?

Постериор

$$P(\theta | D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.
 - При $\tau^2 \rightarrow \infty$ (слабый априор): $\mu_N \rightarrow \hat{\theta}_{\text{MLE}}$.
 - При малом τ^2 : μ_N тяготеет к нулю.
- В целом:

Интуиция: что означают формулы?

Постериор

$$P(\theta | D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.
 - При $\tau^2 \rightarrow \infty$ (слабый априор): $\mu_N \rightarrow \hat{\theta}_{\text{MLE}}$.
 - При малом τ^2 : μ_N тяготеет к нулю.
- В целом:
 - Байесовская регрессия = MLE + априор → неопределённость.

Интуиция: что означают формулы?

Постериор

$$P(\theta | D) = \mathcal{N}(\mu_N, \Sigma_N),$$

где

$$\Sigma_N = \left(\frac{1}{\sigma^2} X^\top X + \frac{1}{\tau^2} I \right)^{-1}, \quad \mu_N = \frac{1}{\sigma^2} \Sigma_N X^\top y.$$

- Σ_N — ковариация постериора.
 - Чем больше данных n , тем больше $X^\top X$, тем меньше дисперсия → растёт уверенность.
 - Малое τ^2 (сильный априор) уменьшает вариацию → сжатие к нулю.
- μ_N — центр постериора.
 - Баланс между MLE и априором.
 - При $\tau^2 \rightarrow \infty$ (слабый априор): $\mu_N \rightarrow \hat{\theta}_{\text{MLE}}$.
 - При малом τ^2 : μ_N тяготеет к нулю.
- В целом:
 - Байесовская регрессия = MLE + априор → неопределённость.
 - Мы получаем распределение параметров, а не одно число.

Определение

Для нового объекта x^* :

$$P(y^* | x^*, D) = \int P(y^* | x^*, \theta, \sigma^2) P(\theta | D) d\theta.$$

- Внутренний интеграл усредняет предсказания по всем возможным θ .
- Благодаря сопряжённости (нормальное \times нормальное) интеграл считается аналитически.

Определение

Для нового объекта x^* :

$$P(y^* | x^*, D) = \int P(y^* | x^*, \theta, \sigma^2) P(\theta | D) d\theta.$$

- Внутренний интеграл усредняет предсказания по всем возможным θ .
- Благодаря сопряжённости (нормальное \times нормальное) интеграл считается аналитически.

Определение

Для нового объекта x^* :

$$P(y^* | x^*, D) = \int P(y^* | x^*, \theta, \sigma^2) P(\theta | D) d\theta.$$

- Внутренний интеграл усредняет предсказания по всем возможным θ .
- Благодаря сопряжённости (нормальное \times нормальное) интеграл считается аналитически.

Результат

$$P(y^* | x^*, D) = \mathcal{N}(x^{*\top} \mu_N, x^{*\top} \Sigma_N x^* + \sigma^2),$$

где μ_N, Σ_N — параметры постериора для θ .

Модель и априор

Линейная регрессия: $y = X\theta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

Гауссовский априор на коэффициенты: $\theta \sim \mathcal{N}(0, \tau^2 I)$

- Постериор по теореме Байеса:

$$P(\theta \mid D) \propto P(y \mid X, \theta) P(\theta) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2\right).$$

Модель и априор

Линейная регрессия: $y = X\theta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

Гауссовский априор на коэффициенты: $\theta \sim \mathcal{N}(0, \tau^2 I)$

- Постериор по теореме Байеса:

$$P(\theta | D) \propto P(y | X, \theta) P(\theta) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2\right).$$

- MAP** = максимум постериора = минимум отрицательного лог-постериора:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \left[\frac{1}{2\sigma^2} \|y - X\theta\|^2 + \frac{1}{2\tau^2} \|\theta\|^2 \right].$$

Модель и априор

Линейная регрессия: $y = X\theta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

Гауссовский априор на коэффициенты: $\theta \sim \mathcal{N}(0, \tau^2 I)$

- Постериор по теореме Байеса:

$$P(\theta | D) \propto P(y | X, \theta) P(\theta) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2\right).$$

- MAP** = максимум постериора = минимум отрицательного лог-постериора:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \left[\frac{1}{2\sigma^2} \|y - X\theta\|^2 + \frac{1}{2\tau^2} \|\theta\|^2 \right].$$

- Умножим на $2\sigma^2$ (не меняет минимума) и положим $\lambda = \sigma^2/\tau^2$:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \|y - X\theta\|^2 + \lambda \|\theta\|^2.$$

Модель и априор

Линейная регрессия: $y = X\theta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

Гауссовский априор на коэффициенты: $\theta \sim \mathcal{N}(0, \tau^2 I)$

- Постериор по теореме Байеса:

$$P(\theta | D) \propto P(y | X, \theta) P(\theta) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\tau^2} \|\theta\|^2\right).$$

- MAP** = максимум постериора = минимум отрицательного лог-постериора:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \left[\frac{1}{2\sigma^2} \|y - X\theta\|^2 + \frac{1}{2\tau^2} \|\theta\|^2 \right].$$

- Умножим на $2\sigma^2$ (не меняет минимума) и положим $\lambda = \sigma^2/\tau^2$:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \|y - X\theta\|^2 + \lambda \|\theta\|^2.$$

- Это в точности задача **Ridge regression** (Тихоновская регуляризация, L_2):

$$\min_{\theta} \|y - X\theta\|^2 + \lambda \|\theta\|^2.$$

Напомним:

MAP для гауссовского априора эквивалентен Ridge:

$$\hat{\theta}_{\text{MAP}} = (X^T X + \lambda I)^{-1} X^T y, \quad \lambda = \frac{\sigma^2}{\tau^2}.$$

- Вопрос: как выбрать λ ?

Напомним:

MAP для гауссовского априора эквивалентен Ridge:

$$\hat{\theta}_{\text{MAP}} = (X^T X + \lambda I)^{-1} X^T y, \quad \lambda = \frac{\sigma^2}{\tau^2}.$$

- Вопрос: как выбрать λ ?
- Частотный подход: кросс-валидация.

Выбор λ через evidence

Напомним:

MAP для гауссовского априора эквивалентен Ridge:

$$\hat{\theta}_{\text{MAP}} = (X^T X + \lambda I)^{-1} X^T y, \quad \lambda = \frac{\sigma^2}{\tau^2}.$$

- Вопрос: как выбрать λ ?
- Частотный подход: кросс-валидация.
- Байесовский подход: **максимизация evidence**.

$$P(y | X, \lambda, \sigma^2) = \int P(y | X, \theta, \sigma^2) P(\theta | \lambda) d\theta.$$

Сопряжённость (нормальное \times нормальное)

Интеграл берётся в замкнутой форме:

$$P(y | X, \lambda, \sigma^2) = \mathcal{N}(y | 0, \sigma^2 I + \tau^2 X X^\top).$$

- В терминах $\lambda = \sigma^2 / \tau^2$:

$$\log P(y | X, \lambda, \sigma^2) = -\frac{1}{2} \left(n \log(2\pi) + \log \det C + y^\top C^{-1} y \right),$$

где $C = \sigma^2 I + \tau^2 X X^\top$.

Сопряжённость (нормальное \times нормальное)

Интеграл берётся в замкнутой форме:

$$P(y | X, \lambda, \sigma^2) = \mathcal{N}(y | 0, \sigma^2 I + \tau^2 X X^\top).$$

- В терминах $\lambda = \sigma^2 / \tau^2$:

$$\log P(y | X, \lambda, \sigma^2) = -\frac{1}{2} \left(n \log(2\pi) + \log \det C + y^\top C^{-1} y \right),$$

где $C = \sigma^2 I + \tau^2 X X^\top$.

- Задача: выбрать $\hat{\lambda} = \arg \max_{\lambda} \log P(y | X, \lambda, \sigma^2)$.

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.
- Большое λ (сильный априор, τ^2 маленькое) \Rightarrow коэффициенты сжаты к нулю, риск недообучения.

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.
- Большое λ (сильный априор, τ^2 маленькое) \Rightarrow коэффициенты сжаты к нулю, риск недообучения.
- Evidence $P(y | X, \lambda)$ балансирует эти два эффекта:

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.
- Большое λ (сильный априор, τ^2 маленькое) \Rightarrow коэффициенты сжаты к нулю, риск недообучения.
- Evidence $P(y | X, \lambda)$ балансирует эти два эффекта:
 - штрафует за избыточную сложность (бритва Оккама),

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.
- Большое λ (сильный априор, τ^2 маленькое) \Rightarrow коэффициенты сжаты к нулю, риск недообучения.
- Evidence $P(y | X, \lambda)$ балансирует эти два эффекта:
 - штрафует за избыточную сложность (бритва Оккама),
 - поощряет модели, которые хорошо объясняют данные.

Интуиция выбора λ через evidence

- Малое λ (слабый априор, τ^2 большое) \Rightarrow модель гибкая, риск переобучения.
- Большое λ (сильный априор, τ^2 маленькое) \Rightarrow коэффициенты сжаты к нулю, риск недообучения.
- Evidence $P(y | X, \lambda)$ балансирует эти два эффекта:
 - штрафует за избыточную сложность (бритва Оккама),
 - поощряет модели, которые хорошо объясняют данные.
- В итоге получаем «оптимальное» λ без кросс-валидации.