

Lezione 3: Richiami di statistica

Nella lezione precedente abbiamo assunto che i parametri $\vec{\theta}$ fossero tutti noti. Questa lezione discute come imporre $\vec{\theta}$ dai dati \mathcal{D} . Come delineato nella prima lezione, lo scopo del gioco si riduce a

$$\hat{\vec{\theta}} = \underset{\vec{\theta}}{\operatorname{argmin}} \mathcal{L}(\vec{\theta})$$

3.1 MAXIMUM LIKELIHOOD ESTIMATION (MLE)

Definiamo la MLE come $\hat{\vec{\theta}}_{\text{mle}} \triangleq \underset{\vec{\theta}}{\operatorname{argmax}} p(\mathcal{D}/\vec{\theta})$. Se i dati sono campioni indipendenti della stessa distribuzione, allora

$$p(\mathcal{D}/\vec{\theta}) = \prod_{n=1}^N p(\vec{y}_n / \vec{x}_n, \vec{\theta})$$

Definiamo la log-likelihood come (negativa \equiv NLL)

$$\ell(\vec{\theta}) \triangleq -\log p(\mathcal{D}/\vec{\theta}) = -\sum_{n=1}^N \log p(\vec{y}_n / \vec{x}_n, \vec{\theta})$$

La MLE è data da

$$\hat{\vec{\theta}}_{\text{mle}} = \underset{\vec{\theta}}{\operatorname{argmin}} -\sum_{n=1}^N \log p(\vec{y}_n / \vec{x}_n, \vec{\theta}) \quad (\text{supervised learning})$$

$$\hat{\vec{\theta}}_{\text{mle}} = \underset{\vec{\theta}}{\operatorname{argmin}} -\sum_{n=1}^N \log p(\vec{y}_n / \vec{\theta}) \quad (\text{unsupervised learning})$$

La MLE può essere giustificata pensando come approssimazione del "posterior" bayesiano, dato un "prior" uniforme. Per esempio:

$$p(\vec{\theta}/\mathcal{D}) = \delta(\vec{\theta} - \hat{\vec{\theta}}_{\text{MAP}}) \quad \text{con } \hat{\vec{\theta}}_{\text{MAP}} \text{ il "posterior"}$$

$$\hat{\vec{\theta}}_{\text{MAP}} = \underset{\vec{\theta}}{\operatorname{argmin}} -\log p(\vec{\theta}/\mathcal{D}) = \underset{\vec{\theta}}{\operatorname{argmin}} -\log p(\mathcal{D}/\vec{\theta}) - \log p(\vec{\theta})$$

$$p(\vec{\theta}) = \mathcal{I} \hat{\vec{\theta}}_{\text{MLE}}$$

Supponiamo che Y sia una variabile casuale distribuita normalmente $Y \sim \mathcal{N}(\mu, \sigma^2)$ e sia \mathcal{D} un dataset i.i.d. punti sono campionati indipendentemente: $\mathcal{D} = \{y_n : n = 1 : N\}$. Allora

$$\begin{aligned} \text{NLL}(\mu, \sigma^2) &= - \sum_{n=1}^{N_D} \ln \left[\left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left(-\frac{1}{2\sigma^2} (y_n - \mu)^2 \right) \right] \\ &= \frac{1}{2\sigma^2} \sum_{n=1}^{N_D} (y_n - \mu)^2 + \frac{N_D}{2} \ln(2\pi\sigma^2) \end{aligned}$$

$$\frac{\partial}{\partial \mu} \text{NLL}(\mu, \sigma^2) = 0 \iff \hat{\mu}_{\text{MLE}} = \frac{1}{N_D} \sum_{n=1}^{N_D} y_n = \bar{y}$$

$$\frac{\partial}{\partial \sigma^2} \text{NLL}(\mu, \sigma^2) = 0 \iff \hat{\sigma}_{\text{MLE}}^2 = \frac{1}{N_D} \sum_{n=1}^{N_D} (y_n - \hat{\mu}_{\text{MLE}})^2 = \underbrace{\frac{1}{N} \sum_{n=1}^{N_D} y_n^2 - \bar{y}^2}_{s^2}$$

Per una distribuzione multivariata:

$$\ell(\bar{\mu}^{-1}, \bar{\Sigma}^{-1}) = \ln p(\mathcal{D} | \bar{\mu}^{-1}, \bar{\Sigma}^{-1}) = \frac{N_D}{2} \ln |\bar{\Lambda}^{-1}| = -\frac{1}{2} \sum_{n=1}^{N_D} (\bar{y}_n^{-1} - \bar{\mu}^{-1})^T \bar{\Lambda} (\bar{y}_n^{-1} - \bar{\mu}^{-1})$$

con $\bar{\Lambda} = \bar{\Sigma}^{-1}$ la PRECISION MATRIX. segue (come sopra)

$$\hat{\mu} = \frac{1}{N_D} \sum_{n=1}^{N_D} \bar{y}_n^{-1} = \bar{y}$$

empirical mean

$$\hat{\Sigma} = \frac{1}{N_D} \sum_{n=1}^{N_D} (\bar{y}_n - \bar{y})(\bar{y}_n - \bar{y})^T$$

empirical covariance matrix

3.2 EMPIRICAL RISK MINIMISATION (ERM)

La MLE può essere generalizzata sostituendo la loss function logaritmica con qualunque altra funzione

$$\mathcal{L}(\bar{\vartheta}) = \frac{1}{N} \sum_{n=1}^N \ell(\bar{y}_n^{-1}, \bar{\vartheta}, \bar{x}_n^{-1})$$

Cio' è noto come EMPIRICAL RISK MINIMISATION (ERM) dal momento che è il loro atteso quando l'aspettazione è presa rispetto alla distribuzione empirica. Per esempio, in un problema di classificazione

$$\ell_{01}(\bar{y}_n^{-1}, \bar{\vartheta}, \bar{x}_n^{-1}) = \begin{cases} 0 & \text{se } \bar{y}_n^{-1} = f(\bar{x}_n^{-1}, \bar{\vartheta}) \\ 1 & \text{se } \bar{y}_n^{-1} \neq f(\bar{x}_n^{-1}, \bar{\vartheta}) \end{cases}$$

con $f(\vec{x}, \vec{\theta})$ un predittore. L'empirical risk diventa

$$\mathcal{L}(\vec{\theta}) = \frac{1}{N} \sum_{n=1}^N l_{0,1}(\vec{y}_n, \vec{\theta}; \vec{x}_n)$$

che è il misclassification rate (ml training set). Per problemi binari possiamo riscrivere il misclassification rate nella forma seguente

$$\tilde{y} \in \{-1, +1\} \text{ true label} \quad \hat{y} \in \{-1, +1\} = f(\vec{x}, \vec{\theta}) \text{ prediction}$$

$$l_{0,1}(\tilde{y}, \hat{y}) = \mathbb{1}(\tilde{y} \neq \hat{y}) = \mathbb{1}(\tilde{y} \hat{y} < 0)$$

Il rischio empirico diventa

$$\mathcal{L}(\vec{\theta}) = \frac{1}{N} \sum_{n=1}^N l_{0,1}(y_n, \hat{y}_n) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(\tilde{y}_n, \hat{y}_n < 0)$$

La funzione $l_{0,1}$ è non-smooth, il che la rende difficile da ottimizzare. In luogo di essa, si può utilizzare una funzione smooth (definita generalmente come limite superiore convessa). Per esempio

$$l(\tilde{y}, \eta) = -\ln p(\tilde{y}|\eta) = \ln(1 + e^{-\tilde{y}\eta}) \quad p(\tilde{y}|\vec{x}, \vec{\theta}) = \sigma(\tilde{y}\eta) = \frac{1}{1 + e^{-\tilde{y}\eta}}$$

Altri esempi sono l'average loss e l'exp loss.

3 ALTRI METODI DI MINIMIZZAZIONE

3.1 METODO DEI MOMENTI

Valore il calcolo di $\nabla_{\vec{\theta}} \text{MLL}(\vec{\theta}) = 0$ è difficile. Il metodo dei momenti consiste nell'equagliare i momenti teorici della distribuzione ai momenti empirici.

$$\begin{array}{ll} \text{MOMENTI TEORICI:} & \mu_k = \mathbb{E}[Y^k] \\ \text{MOMENTI EMPIRICI} & \hat{\mu}_k = \frac{1}{N} \sum_{n=1}^N y_n^k \end{array} \left\{ \begin{array}{l} \mu_k = \hat{\mu}_k \text{ } k \text{ equazioni} \\ \text{ho un set di} \end{array} \right.$$

Per la gaussiana:

$$\mu_1 = \mu = \bar{y}$$

$$\mu_2 = \sigma^2 + \mu^2 = s^2$$

in questo caso $\text{MOM} \equiv \text{MLE}$. Non sempre è così (p. es. distr. uniforme)

3.3.2. ONLINE RECURSIVE ESTIMATION

4

Se tutto il dataset \mathcal{D} è noto e disponibile prima che il learning venga inizializzato, si dice che si fa "batch learning". In alcuni casi, però, il dataset è disponibile a blocchi.

$\hat{\theta}_{t-1}$ è la predizione dato $\mathcal{D}_{1:t-1}$

dobbiamo trovare $\hat{\theta}_t = f(\hat{\theta}_{t-1}, y_t)$ con un update ricorsivo

Per una gaussiana univariata

$$\hat{\mu}_t = \frac{1}{t} \sum_{n=1}^t y_n \longrightarrow \hat{\mu}_t = \frac{1}{t} \sum_{n=1}^t y_n = \frac{1}{t} ((t-1)\hat{\mu}_{t-1} + y_t)$$

$$= \hat{\mu}_{t-1} + \frac{1}{t} (y_t - \hat{\mu}_{t-1}) \quad \text{MOVING AVERAGE}$$

La moving average può eventualmente essere pesata se la distribuzione cambia cont.

3.4 REGOLARIZZAZIONE

Un problema della MLE e della ERM è che i parametri tenderanno ad essere determinati minimizzando il loss sul training set, ma non necessariamente sui dati futuri. Il problema si risolve con la regolarizzazione

$$\mathcal{L}(\bar{\theta}, \lambda) = \left[\frac{1}{N} \sum_{n=1}^N \ell(\bar{y}_n, \bar{\theta}; \bar{x}_n) + \lambda C(\bar{\theta}) \right]$$

\uparrow \uparrow
 regularisation penalty
 parameter function
 p.es. $C(\bar{\theta}) = -\ln p(\bar{\theta})$

Per $\lambda = 1$ e riscalandolo $p(\bar{\theta})$, si ha la NLL

$$NLL(\bar{\theta}, \lambda) = - \left[\sum_{n=1}^N \ln p(\bar{y}_n / \bar{x}_n, \bar{\theta}) + \ln p(\bar{\theta}) \right] = - \left[\ln p(\mathcal{D} / \bar{\theta}) + \ln p(\bar{\theta}) \right]$$

che implica

$$\hat{\theta} = \underset{\bar{\theta}}{\operatorname{argmax}} \ln p(\bar{\theta} / \mathcal{D}) = \underset{\bar{\theta}}{\operatorname{argmax}} \left[\ln p(\mathcal{D} / \bar{\theta}) + \ln p(\bar{\theta}) - \text{const} \right]$$

Questa si chiama MAP (MAXIMUM A POSTERIOR) estimation.

Come scegliere il valore di λ ? Un valore (troppo) piccolo significa minimizzare il rischio empirico, un valore (troppo) grande significa essere troppo vicini al prior (overfitting vs underfitting).
 Dividiamo il dataset in due classi: training e validation (80/20%)
 Si fitta il modello su \mathcal{D}_{train} \forall setting λ e si valuta la performance R_{valid} . Si prende quindi il valore di λ associato alle performance migliori. Definiamo l'empirical risk regolarizzato

$$R_{\lambda}(\bar{\theta}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\bar{x}, \bar{y}) \in \mathcal{D}} \ell(\bar{y}, f(\bar{x}, \bar{\theta})) + \lambda C(\bar{\theta})$$

$$\hat{\theta}_{\lambda}(\mathcal{D}_{train}) = \underset{\bar{\theta}}{\operatorname{argmin}} R_{\lambda}(\bar{\theta}, \mathcal{D}_{train}) \quad \forall \lambda$$

$$R_{\lambda}^{val} \triangleq R_0(\hat{\theta}_{\lambda}(\mathcal{D}_{train}), \mathcal{D}_{valid}) \quad \text{si calcola il VALIDATION RISK}$$

$$\lambda^* = \underset{\lambda \in \mathcal{S}}{\operatorname{argmin}} R_{\lambda}^{val} \quad \text{Dopo aver preso } \lambda^* \text{ si rifitte il modello su } \mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{valid}$$

$$\hat{\theta}^* = \underset{\bar{\theta}}{\operatorname{argmin}} R_{\lambda^*}(\bar{\theta}, \mathcal{D})$$

Se il dataset è piccolo, eliminare il 20% dei dati dal training può essere dannoso. Si usa allora il CROSS-VALIDATION in cui il dataset di training è diviso in K folds. Per ogni fold, il modello viene allenato su tutti i folds, tranne il K -esimo, che viene usato come test fold.

$$R_{\lambda}^{cv} \triangleq \frac{1}{K} \sum_{k=1}^K R_0(\hat{\theta}_{\lambda}(\mathcal{D} - \mathcal{D}_k), \mathcal{D}_k)$$

6 STATISTICA BAYESIANA

Abbiamo visto vari metodi per determinare i parametri dai dati, ma non abbiamo detto nulla sulle loro incertezze. Il TB ci aiuta

$$p(\bar{\theta}' / \mathcal{D}) = \frac{p(\bar{\theta}') p(\mathcal{D} / \bar{\theta}')}{p(\mathcal{D})} = \frac{p(\bar{\theta}') p(\mathcal{D} / \bar{\theta}')}{\int p(\bar{\theta}') p(\mathcal{D} / \bar{\theta}') d\bar{\theta}'}$$

$$\mathcal{D} = \{(\bar{x}_n, \bar{y}_n)\} \quad \text{SUPERVISED LEARNING}$$

$$n = 1, \dots, N$$

$$\mathcal{D} = \{(\bar{y}_n)\} \quad \text{UNSUPERVISED LEARNING}$$

Una volta determinato il posteriore sui parametri, possiamo calcolare il posteriore delle distribuzioni predittive marginalizzando su $\bar{\theta}$

$$p(\bar{y} | \bar{x}, \mathcal{D}) = \int p(\bar{y} | \bar{x}, \bar{\theta}) p(\bar{\theta} | \mathcal{D}) d\bar{\theta} \quad \text{Bayes model averaging (BMA)}$$

Consideriamo una distribuzione gaussiana di cui sia nota la varianza. Nel caso univariato, la likelihood per μ ha la forma

$$p(\mathcal{D} | \mu) \propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^{N_D} (y_n - \mu)^2 \right\}$$

Ma può mostrare che il prior coniugato è un'altra gaussiana $\mathcal{N}(\mu | \tilde{m}, \tilde{\tau}^2)$. Usando il teorema di Bayes, troviamo che il posteriore è

$$p(\mu | \mathcal{D}, \sigma^2) = \mathcal{N}(\mu | \hat{m}, \hat{\tau}^2)$$

$$\hat{\tau}^2 = \frac{1}{\frac{N}{\sigma^2} + \frac{1}{\tilde{\tau}^2}} = \frac{\sigma^2 \tilde{\tau}^2}{N\tilde{\tau}^2 + \sigma^2} \quad \hat{m} = \hat{\tau}^2 \left(\frac{\tilde{m}}{\tilde{\tau}^2} + \frac{N\bar{y}}{\sigma^2} \right) = \frac{\sigma^2}{N\tilde{\tau}^2 + \sigma^2} \tilde{m} + \frac{N\tilde{\tau}^2}{N\tilde{\tau}^2 + \sigma^2} \bar{y}$$

con $\bar{y} \triangleq \frac{1}{N} \sum_{n=1}^N y_n$ è la media empirica.

Definendo $K = 1/\sigma^2$ e $\tilde{\lambda} = 1/\tilde{\tau}^2$ ottengo

$$\hat{\lambda} = \tilde{\lambda} + NK$$

$$\hat{m} = \frac{NK\bar{y} + \tilde{\lambda}\tilde{m}}{\hat{\lambda}} = \frac{NK}{NK + \tilde{\lambda}} \bar{y} + \frac{\tilde{\lambda}}{NK + \tilde{\lambda}} \tilde{m}$$

La precisione del "posterior" è la precisione del "prior" più N volte la precisione K . La media del "posterior" è una combinazione convessa delle medie empirica \bar{y} e delle

media del prior \bar{m} . Consideriamo ora il posteriore dopo aver visto un singolo data point y (quindi $N=1$). Allora

$$\hat{m} = \underbrace{\frac{1}{\lambda} \bar{m}}_{\text{convex combination of prior and data}} + \underbrace{\frac{K}{\lambda} \bar{y}}_{\text{prior mean adjusted to data}} = \bar{y} - \underbrace{\frac{1}{\lambda} (\bar{y} - \bar{m})}_{\text{data adjusted to the prior mean}}$$

POSTERIOR MEAN

$$se(\mu) = \sqrt{V[\mu|\mathcal{D}]} \quad \text{STANDARD ERROR}$$

Se usi un "uninformative prior" per μ ponendo $\lambda=0$, allora $\hat{m} = \bar{y}$.
 Supponendo di approssimare $\sigma^2 \sim s^2 \approx \frac{1}{N} \sum_{n=1}^N (y_n - \bar{y})^2$ segue che $\lambda = NK = N/s^2 \Rightarrow se(\mu) = \frac{1}{\sqrt{\lambda}} = \frac{s}{\sqrt{N}}$.

L'incertezza su μ si riduce ad un rateo di $1/\sqrt{N}$.

Osservazioni:

- 1) Quando non abbiamo informazioni sul prior, è desiderabile usare un prior "uninformative". Per esempio un flat prior $p(\mu)=1$.
- 2) Qualunque modello Bayesiano richiede che venga specificato un prior $p(\bar{\theta})$ per i parametri. I parametri del prior vengono chiamati iperparametri (e sono denotati con Φ). Se non sono noti, possiamo metterci sopra un prior (multi-level model $\Phi \rightarrow \bar{\theta} \rightarrow \mathcal{D}$)
 La joint distribution è $p(\bar{\Phi}, \bar{\theta}, \mathcal{D}) = p(\Phi) p(\bar{\theta}|\Phi) p(\mathcal{D}|\bar{\theta})$.

METHOD	DEFINITION
MLE	$\hat{\bar{\theta}} = \operatorname{argmax}_{\bar{\theta}} p(\mathcal{D} \bar{\theta})$
MAP	$\hat{\bar{\theta}} = \operatorname{argmax}_{\bar{\theta}} p(\mathcal{D} \bar{\theta}) p(\bar{\theta} \Phi)$
FULL BAYES	$p(\bar{\theta}, \Phi \mathcal{D}) \propto p(\mathcal{D} \bar{\theta}) p(\bar{\theta} \Phi) p(\Phi)$

Una distribuzione posteriore è un oggetto multidimensionale difficile da visualizzare e da trattare. Può essere utile quindi calcolare stimatori puntuali (come il posterior mean e mode) e calcolare un intervallo di credibilità che quantifica l'incertezza associata a quelle stime. L'intervallo di credibilità $100(1-\alpha)\%$ è la regione $C = (l, u)$ che contiene $1-\alpha$ delle probabilità posteriori.

$$C_\alpha(\mathcal{D}) = (l, u) : P(l \leq \vartheta \leq u | \mathcal{D}) = 1 - \alpha$$

l : lower
 u : upper

Se il "posterior" ha una PDF nota, allora $l = F^{-1}(\alpha/2)$ e $u = F^{-1}(1-\alpha/2)$ dove F è la CDF del "posterior". Esempio: gaussiano.

Osservazione

L'intervallo di credibilità è un concetto Bayesiano. L'intervallo di confidenza è un concetto frequentista. Si definisce intervallo di confidenza $100(1-\alpha)\%$ per le stime di un parametro ϑ l'intervallo $I(\tilde{\mathcal{D}}) = (l(\tilde{\mathcal{D}}), u(\tilde{\mathcal{D}}))$ ottenuto da un data set $\tilde{\mathcal{D}}$ tale che

$$\Pr(\vartheta \in I(\tilde{\mathcal{D}}) | \tilde{\mathcal{D}} \sim \vartheta) = 1 - \alpha$$

Un CI al 95% per un parametro ϑ non significa che il parametro stia verosimilmente dentro l'intervallo il 95% delle volte dati i dati osservati. CI (frequentista): ϑ è una costante fissa non nota e i dati sono aleatori; CI (Bayesiano): i dati sono fissi perché noti, mentre ϑ è ignota.

Esempio

Supponiamo di generare due interi $\mathcal{D} = (y_1, y_2)$ da

$$p(y|\vartheta) = \begin{cases} 0,5 & \text{se } y = \vartheta \\ 0,5 & \text{se } y = \vartheta + 1 \\ 0 & \text{altrimenti} \end{cases}$$

Se $\vartheta = 33$, ci aspettiamo i seguenti risultati:

$$(33, 33), (33, 40), (40, 33), (40, 40)$$

Sia $m = \min(y_1, y_2)$ e l'intervallo $[l(\vartheta), u(\vartheta)] = [m, m]$. Segue

$$[33, 33], [33, 33], [33, 33], [40, 40] \Rightarrow 75\% \text{ CI for } \vartheta = 33.$$

Ma se osserviamo $\mathcal{D} = (33, 40)$, allora $p(\vartheta = 33 | \mathcal{D}) = 1$ (ma ci è solo 75%)

Il CI fallisce per esperimenti NON ripetibili.

3.8 BIAS-VARIANCE TRADEOFF

Sia $\hat{\vartheta}$ lo stimatore statistico e ϑ^* l'estimando. Nel formalismo frequentista, i dati sono variabili casuali, campionati da una distribuzione $p^*(\mathcal{D})$ che induce una distribuzione sull'estimando $p^*(\hat{\vartheta}(\mathcal{D}))$

Il BIAS di uno stimatore è definito come

$$\text{bias}(\hat{\vartheta}(\cdot)) \triangleq \mathbb{E}[\hat{\vartheta}(\mathcal{D})] - \vartheta^* \quad (\vartheta^* \text{ è il valore vero})$$

Se il BIAS è nullo, lo stimatore viene detto UNBIASED. Per una popolazione

$$\text{bias}(\hat{\mu}) = \mathbb{E}[\bar{x}] - \mu = \mathbb{E}\left[\frac{1}{N_D} \sum_{n=1}^{N_D} x_n\right] - \mu = \frac{N_D \mu}{N_D} - \mu = 0 \quad \text{UNBIASED}$$

$$\text{bias}(\hat{\sigma}^2) = \mathbb{E}[\mathbf{s}^2] - \sigma^2 = \frac{N_D - 1}{N_D} \sigma^2 - \sigma^2 \neq 0 \quad \text{BIASED}$$

$$\Rightarrow \hat{\sigma}^2 \rightarrow \hat{\sigma}_{\text{mlt}}^2 = \frac{1}{N_D - 1} \sum_{n=1}^{N_D} (x_n - \bar{x})^2 = \frac{N_D}{N_D - 1} \sigma_{\text{MLE}}^2$$

La VARIANZA di uno stimatore è

$$V[\hat{\vartheta}] \triangleq \mathbb{E}[\hat{\vartheta}^2] - (\mathbb{E}[\hat{\vartheta}])^2$$

Idealmente vogliamo che V sia minima \forall stimatore. Il teorema di Cramer-Rao fornisce un limite inferiore alle varianze

per uno stimatore UNBIASED ϑ^* :

$$V[\hat{\vartheta}] \geq \frac{1}{NF(\vartheta^*)} \quad \text{dove } F(\vartheta^*) \text{ è la Fisher information matrix}$$

Si può dimostrare che la MLE raggiunge il bound di Cramer-Rao.

Lo scopo dell'inferenza statistica è minimizzare l'errore quadratico medio (MSE). Sia $\hat{\vartheta} = \hat{\vartheta}(\mathcal{D})$ l'estimatore e $\bar{\vartheta} = E[\hat{\vartheta}]$ il valore di aspettazione corrispondente (tutti sotto $p(\mathcal{D}|\vartheta^*)$). Allora

$$\begin{aligned} E[(\hat{\vartheta} - \vartheta^*)^2] &= E[(\hat{\vartheta} - \bar{\vartheta}) + (\bar{\vartheta} - \vartheta^*)]^2 \\ &= E[(\hat{\vartheta} - \bar{\vartheta})^2] + 2(\bar{\vartheta} - \vartheta^*) E[\hat{\vartheta} - \bar{\vartheta}] + (\bar{\vartheta} - \vartheta^*)^2 \\ &= E[(\hat{\vartheta} - \bar{\vartheta})^2] + (\bar{\vartheta} - \vartheta^*)^2 = V[\hat{\vartheta}] + \text{bias}^2(\hat{\vartheta}) \end{aligned}$$

cioè $MSE = \text{variance} + \text{bias}^2$. BIAS-VARIANCE TRADE OFF

Esempio

Supponiamo di stimare le medie di una distribuzione gaussiana per $\bar{x} = (x_1, \dots, x_N)$ assumendo che i dati siano campionati da $x_n \sim \mathcal{N}(\vartheta^* = \mu, \sigma^2)$. Uno stimatore ovvio è dato dalla MLE che ha bias 0 e varianza $V[\bar{x}|\vartheta^*] = \frac{\sigma^2}{N}$. Ma possiamo avere uno stimatore MAP sotto un prior gaussiano. In tal caso

$$\tilde{x} = \frac{N}{N+K_0} \bar{x} + \frac{K_0}{N+K_0} \vartheta_0 = w \bar{x} + (1-w) \vartheta_0 \quad \begin{array}{l} \text{con un prior} \\ \mathcal{N}(\vartheta_0, \frac{\sigma^2}{K_0}) \end{array}$$

$$\begin{aligned} E[\tilde{x}] - \vartheta^* &= w \vartheta^* + (1-w) \vartheta_0 - \vartheta^* = (1-w)(\vartheta_0 - \vartheta^*) \\ V[\tilde{x}] &= w^2 \frac{\sigma^2}{N} < V_{MLE}[\bar{x}] \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \begin{array}{l} \text{BIASED} \\ 0 < w < 1 \end{array}$$