# 1 Fundamentals

- **Normal**: $\dfrac{\exp(-\frac{1}{2}(\mathbf{x}-\mu)^T\Sigma^{-1}(\mathbf{x}-\mu))}{\sqrt{(2\pi)^k\det(\Sigma)}}$
- **Beta**: $\text{Beta}(\theta;\alpha,\beta)\propto\theta^{\alpha-1}(1-\theta)^{\beta-1}$
- **Laplace**: $\frac{1}{2l}\exp(-\frac{|x-\mu|}{l})$
- Gaussian CDF has no closed-form; $O(n^2)$ params.
- $\mathbb{E}[\mathbf{AX}+\mathbf{b}]=\mathbf{A}\mathbb{E}[\mathbf{X}]+\mathbf{b}$; $\mathbb{E}[\mathbf{X}+\mathbf{Y}]=\mathbb{E}[\mathbf{X}]+\mathbb{E}[\mathbf{Y}]$
- $\mathbb{E}[\mathbf{XY}^\top]=\mathbb{E}[\mathbf{X}]\cdot\mathbb{E}[\mathbf{Y}]^\top$ (if $\mathbf{X},\mathbf{Y}$ indep.)
- LOTUS: $\mathbb{E}[\mathbf{g}(\mathbf{X})]=\int_{\mathbf{X}(\Omega)}\mathbf{g}(\mathbf{x})\cdot p(\mathbf{x})d\mathbf{x}$ (if $\mathbf{g}$ nice and $\mathbf{X}$ cont.)
- Tower rule: $\mathbb{E}_\mathbf{Y}[\mathbb{E}_\mathbf{X}[\mathbf{X}|\mathbf{Y}]]=\mathbb{E}[\mathbf{X}]$
- $\text{Var}[\mathbf{X}]=\mathbb{E}[(\mathbf{X}-\mathbb{E}[\mathbf{X}])(\mathbf{X}-\mathbb{E}[\mathbf{X}])^\top]$
  $=\mathbb{E}[\mathbf{XX}^\top]-\mathbb{E}[\mathbf{X}]\cdot\mathbb{E}[\mathbf{X}]^\top=\text{Cov}[\mathbf{X},\mathbf{X}]$
- $\text{Var}[\mathbf{AX}+\mathbf{b}]=\mathbf{A}\text{Var}[\mathbf{X}]\mathbf{A}^\top$
- $\text{Var}[\mathbf{X}+\mathbf{Y}]=\text{Var}[\mathbf{X}]+\text{Var}[\mathbf{Y}]+2\text{Cov}[\mathbf{X},\mathbf{Y}]$
- $\text{Var}[\mathbf{X}+\mathbf{Y}]=\text{Var}[\mathbf{X}]+\text{Var}[\mathbf{Y}]$ (if $\mathbf{X},\mathbf{Y}$ indep.)
- Law of total variance, LOTV:
  $\text{Var}[\mathbf{X}]=\mathbb{E}_\mathbf{Y}[\text{Var}_\mathbf{X}[\mathbf{X}|\mathbf{Y}]]+\text{Var}_\mathbf{Y}[\mathbb{E}_\mathbf{X}[\mathbf{X}|\mathbf{Y}]]$
- $\text{Cov}[\mathbf{X},\mathbf{Y}]\doteq\mathbb{E}[(\mathbf{X}-\mathbb{E}[\mathbf{X}])(\mathbf{Y}-\mathbb{E}[\mathbf{Y}])^\top]$
  $=\mathbb{E}[\mathbf{XY}^\top]-\mathbb{E}[\mathbf{X}]\cdot\mathbb{E}[\mathbf{Y}]^\top$
- $\text{Cov}[\mathbf{X},\mathbf{Y}]=\text{Cov}[\mathbf{Y},\mathbf{X}]$; $\text{Cov}[\mathbf{X},\mathbf{Y}]\geq\mathbf{0}$
- $\text{Cov}[\mathbf{AX}+\mathbf{c},\mathbf{BY}+\mathbf{d}]=\mathbf{A}\text{Cov}[\mathbf{X},\mathbf{Y}]\mathbf{B}^\top$
- *Correlation* is normalized covariance:
  $\text{Cor}[\mathbf{X},\mathbf{Y}](i,j)\doteq\frac{\text{Cov}[X_i,Y_j]}{\sqrt{\text{Var}[X_i]\text{Var}[Y_j]}}\in[-1,1]$
- *Uncorrelated* iff $\text{Cov}[\mathbf{X},\mathbf{Y}]=\mathbf{0}$.
- **Change of variables**: Let $\mathbf{g}$ be diff. and inv. Then for
  $\mathbf{Y}=\mathbf{g}(\mathbf{X})$: $p_\mathbf{Y}(\mathbf{y})=p_\mathbf{X}(\mathbf{g}^{-1}(\mathbf{y}))\cdot|\det(\mathbf{Dg}^{-1}(\mathbf{y}))|$
  where $\mathbf{Dg}^{-1}(\mathbf{y})$ is the Jacobian of $\mathbf{g}^{-1}$ at $\mathbf{y}$.

**Bayes' rule**: $p(\mathbf{x}|\mathbf{y})=\frac{p(\mathbf{y}|\mathbf{x})\cdot p(\mathbf{x})}{p(\mathbf{y})}$

- If prior $p(\mathbf{x})$ and posterior $p(\mathbf{x}|\mathbf{y})$ from same fam. of distr., prior is **conjugate prior** to likelihood $p(\mathbf{y}|\mathbf{x})$.
- Beta distr. is a conjugate prior to binomial likelihood.
- Under some conditions, **Gaussian is self-conjugate** (Gaussian prior and likelihood $\to$ posterior Gaussian).

- Choosing non-informative prior in absence of evidence is **principle of indifference/insufficient reason**.
- **Improper prior**: not required that prior is a valid distr. (i.e., integrates to 1). Can still derive meaning.
- **Max. entropy principle**: choose prior s.t. one that makes the least "additional assumptions", i.e., prior least "informative".

**Gaussian properties**
- **Gaussians have max. entropy among all distr.** with known mean and var.: $1/2\cdot\log((2\pi e)^d\det(\Sigma))$
- Jointly Gaussian random vectors, $\mathbf{X}$ and $\mathbf{Y}$, are independent iff $\mathbf{X}$ and $\mathbf{Y}$ are uncorrelated.
- Closed under marginalization and conditioning.

Let $\mathbf{X}$ be Gaussian and index sets $A,B\subseteq[n]$.
For any **marginal distr.** $\mathbf{X}_A\sim\mathcal{N}(\mu_A,\Sigma_{AA})$ and for any **conditional distr.**:
$\mathbf{X}_A|\mathbf{X}_B=\mathbf{x}_B\sim\mathcal{N}(\mu_{A|B},\Sigma_{A|B})$ where
$\mu_{A|B}\doteq\mu_A+\Sigma_{AB}\Sigma_{BB}^{-1}(\mathbf{x}_B-\mu_B)$
$\Sigma_{A|B}\doteq\Sigma_{AA}-\Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$
Observe that the variance can only shrink.
- Additive and closed under affine transformations.
- $\mathbf{M}\cdot\mathcal{N}(\mu,\Sigma)=\mathcal{N}(\mathbf{M}\mu,\mathbf{M}^\top\Sigma\mathbf{M})$
- $\mathcal{N}(\mu_A,\Sigma_A)+\mathcal{N}(\mu_B,\Sigma_B)=\mathcal{N}(\mu_A+\mu_B,\Sigma_A+\Sigma_B)$
- $\mathcal{N}(\mu_A,\Sigma_A)\cdot\mathcal{N}(\mu_B,\Sigma_B)\propto\mathcal{N}(\cdot,\cdot)$

**Maximum likelihood estimate (MLE)**:
$\hat{\theta}_{\text{MLE}}\doteq\arg\max_\theta p(y_{1:n}|\mathbf{x}_{1:n},\theta)=\arg\max\sum_{i=1}^n\log p(y_i|\mathbf{x}_i,\theta)$
- **Consistent** if: $\hat{\theta}_{\text{MLE}}\xrightarrow{\mathbb{P}}\theta^\star$ as $n\to\infty$.
- **Asymptotically normal** if $\hat{\theta}_{\text{MLE}}\xrightarrow{\mathcal{D}}\mathcal{N}(\theta^\star,\mathbf{S}_n)$ as $n\to\infty$ where $\mathbf{S}_n$ is asymptotic covar. of MLE.

---

- MLE is **asymptotically efficient** (there exists no other consistent estimator with a "smaller" asymptotic var.).
- For the finite sample regime, the MLE need not be unbiased, and it is susceptible to overfitting to the (finite) training data.

**Maximum a posterior (MAP) estimate**:
$\hat{\theta}_{\text{MAP}}\doteq\arg\min_{\theta\in\Theta}\underbrace{-\log p(\theta)}_{\text{regularization}}+\underbrace{\ell_{\text{nll}}(\theta;\mathcal{D}_n)}_{\text{quality of fit}}$
The **log-prior** $\log p(\theta)$ acts as a regularizer. Common:
- $p(\theta)=\mathcal{N}(\theta;\mathbf{0},\lambda\mathbf{I})$ gives $-\log p(\theta)=\frac{\lambda}{2}\|\theta\|_2^2+$const
- $p(\theta)=\text{Laplace}(\theta;\mathbf{0},\lambda)$ gives $-\log p(\theta)=\lambda\|\theta\|_1+$const
- Uniform prior gives const (no reg., MAP is MLE)

# 2 Bayesian Linear Regression (BLR)

$\hat{\mathbf{w}}_{\text{ls}}=(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$;    $\hat{\mathbf{w}}_{\text{ridge}}=(\mathbf{X}^\top\mathbf{X}+\lambda\mathbf{I})^{-1}\mathbf{X}^\top\mathbf{y}$

**Gaussian prior on weights** $\mathbf{w}\sim\mathcal{N}(\mathbf{0},\sigma_p^2\mathbf{I})$:
- Yields Gaussian posterior $\mathbf{w}|\mathbf{x}_{1:n},y_{1:n}\sim\mathcal{N}(\mu,\Sigma)$, as
  $\log p(\mathbf{w}|\mathbf{x}_{1:n},y_{1:n})=-\frac{1}{2}[\mathbf{w}^\top\Sigma^{-1}\mathbf{w}-2\mu]+$const, w.
  $\Sigma\doteq(\sigma_n^{-2}\mathbf{X}^\top\mathbf{X}+\sigma_p^{-2}\mathbf{I})^{-1}$ and $\mu\doteq\sigma_n^{-2}\Sigma\mathbf{X}^\top\mathbf{y}$.
- MAP is *identical to ridge regression* with $\lambda\doteq\sigma_n^2/\sigma_p^2$.
- **Bayesian inference**: Distr. for a test point $\mathbf{x}^\star$ is:
  $y^\star|\mathbf{x}^\star,\mathbf{x}_{1:n},y_{1:n}\sim\mathcal{N}(\mu^\top\mathbf{x}^\star,\mathbf{x}^{\star\top}\Sigma\mathbf{x}^\star+\sigma_n^2)$

**Laplace prior on weights** $\mathbf{w}\sim\text{Laplace}(\mathbf{0},h)$:
- MAP is *identical to lasso regression* with $\lambda\doteq\sigma_n^2/h$.

**Heteroscedastic noise** $\epsilon_i$ may depend on $\mathbf{x}_i$, while **Homoscedastic** may not.

$\text{Var}[y^\star|\mathbf{x}^\star]=\underbrace{\mathbb{E}_\theta[\text{Var}_{y^\star}[y^\star|\mathbf{x}^\star,\theta]]}_{\text{aleatoric uncertainty}}+\underbrace{\text{Var}_\theta[\mathbb{E}_{y^\star}[y^\star|\mathbf{x}^\star,\theta]]}_{\text{epistemic uncertainty}}$

**Aleatoric:** noise in data; **Epistemic:** noise in model.
Applying linear reg. to non-linear fns.: use non-linear transformation $\phi$ to $\mathbf{X}$. Define $\mathbf{\Phi}=\phi(\mathbf{X})$. With Gaussian prior and $\mathbf{K}=\sigma_p^2\mathbf{\Phi}\mathbf{\Phi}^\top$:
$\mathbf{f}|\mathbf{X}\sim\mathcal{N}(\mathbf{\Phi}\mathbb{E}[\mathbf{w}],\mathbf{\Phi}\text{Var}[\mathbf{w}]\mathbf{\Phi}^\top)=\mathcal{N}(\mathbf{0},\mathbf{K})$.

**Kernel**: $k(\mathbf{x},\mathbf{x}')\doteq\sigma_p^2\phi(\mathbf{x})^\top\phi(\mathbf{x}')=\text{Cov}[f(\mathbf{x}),f(\mathbf{x}')]$
- Choice of kernel implicitly determines the function class that $\mathbf{f}$ is sampled from, which encodes our prior beliefs.
- Kernel matrix has shape $n\times n$ (input space dimension) instead of $e\times e$ (feature space dimension).

For **inference**, define $\bar{\mathbf{\Phi}}\doteq\begin{bmatrix}\mathbf{\Phi}\\\phi(\mathbf{x}^\star)^\top\end{bmatrix}$, $\bar{\mathbf{y}}\doteq\begin{bmatrix}\mathbf{y}\\y^\star\end{bmatrix}$, $\bar{\mathbf{f}}\doteq\begin{bmatrix}\mathbf{f}\\f^\star\end{bmatrix}$.
For $\bar{\mathbf{f}}=\bar{\mathbf{\Phi}}w$ we have: $\tilde{\mathbf{y}}|\mathbf{X},\mathbf{x}^\star\sim\mathcal{N}(\mathbf{0},\bar{\mathbf{K}}+\sigma_n^2\mathbf{I})$
- **Linear**: $k(\mathbf{x},\mathbf{x}')=l\mathbf{x}^\top\mathbf{x}'$
- **RBF/Gaussian**: $k(\mathbf{x},\mathbf{x}')=\exp(-\frac{(\mathbf{x}-\mathbf{x}')^2}{2l^2})$ (larger $l$ gives smoother fns.; cannot model under weight-space view of BLR; feature space are poly. of infinite degree)
- **Polynomial** $k(\mathbf{x},\mathbf{x}')=(1+\mathbf{x}^\top\mathbf{x}')^d$ (feature space is poly. of degree $d$)
- **Laplacian**: $k(\mathbf{x},\mathbf{x}')=\exp(-\alpha\|\mathbf{x}-\mathbf{x}'\|)$
- **Matérn**: (For $v=\frac{1}{2}$ Laplace, for $v\to\infty$ RBF)

**Properties of kernels**
- Symmetric: $k(\mathbf{x},\mathbf{x}')=k(\mathbf{x}',\mathbf{x})$    • $\mathbf{K}_{AA}$ is p.s.d.
- **Stationary** if there exists $\tilde{k}$ s.t. $\tilde{k}(\mathbf{x}-\mathbf{x}')=k(\mathbf{x},\mathbf{x}')$ (only relative location of points matters)
- **Isotropic** if there exists $\tilde{k}$ s.t. $\tilde{k}(\|\mathbf{x}-\mathbf{x}'\|_2)=k(\mathbf{x},\mathbf{x}')$ (only distance between points matters)

**Composition of kernels**
- Addition: $k(\mathbf{x},\mathbf{x}')=k_1(\mathbf{x},\mathbf{x}')+k_2(\mathbf{x},\mathbf{x}')$ (OR)
- Multiplication: $k(\mathbf{x},\mathbf{x}')=k_1(\mathbf{x},\mathbf{x}')\cdot k_2(\mathbf{x},\mathbf{x}')$ (AND)
- Mult. with const.: $k(\mathbf{x},\mathbf{x}')\doteq c\cdot k_1(\mathbf{x},\mathbf{x}')$ for any $c\geq0$
- Composition. with poly.: $k(\mathbf{x},\mathbf{x}')\doteq f(k_1(\mathbf{x},\mathbf{x}'))$ for any poly. $f$ with positive coefficients.
- Composition. with exp.: $k(\mathbf{x},\mathbf{x}')\doteq\exp(k_1(\mathbf{x},\mathbf{x}'))$

---

**Efficient online BLR** ($O(d^2)$ instead of $O(d^3)$):
- $\mathbf{X}^{(t+1)^\top}\mathbf{X}^{(t+1)}=\mathbf{X}^{(t)^\top}\mathbf{X}^{(t)}+\mathbf{x}^{(t)^\top}\mathbf{x}^{(t)}\in\mathbb{R}^{d\times d}$
- $\mathbf{X}^{(t+1)^\top}\mathbf{y}^{(t+1)}=\mathbf{X}^{(t)^\top}\mathbf{y}^{(t)}+\mathbf{x}^{(t)^\top}\mathbf{y}^{(t)}\in\mathbb{R}^d$
- Since $\mathbf{X}^\top\mathbf{X}=\sum_{i=1}^t\mathbf{x}_i\mathbf{x}_i^\top$ and $\mathbf{X}^\top\mathbf{y}=\sum_{i=1}^t y_i x_i$

**Logistic BLR**:
$\hat{\mathbf{w}}_{\text{MAP}}\doteq\arg\min_{\mathbf{w}\in\mathbb{R}^d}\frac{1}{2\sigma_p^2}\|\mathbf{w}\|_2^2+\sum_{i=1}^n\log(1+\exp(-y_i\mathbf{w}^\top\mathbf{x}_i))$
- For $\lambda=1/(2\sigma_p^2)$ this is equiv. to standard logistic reg. where $\ell_{\log}(\mathbf{w}^\top\mathbf{x};y)\doteq\log(1+\exp(-y\mathbf{w}^\top\mathbf{x}))$ and $\nabla_\mathbf{w}\ell_{\log}(\mathbf{w}^\top\mathbf{x};y)=-y\mathbf{x}\cdot\sigma(-y\mathbf{w}^\top\mathbf{x})$.
- Post. not Gaussian or closed, but log. density is convex

# 3 Gaussian Processes (GPs)

**Mean fn.** $\mu:\mathcal{X}\to\mathbb{R}$; **Covar. fn.** $k:\mathcal{X}\times\mathcal{X}\to\mathbb{R}$;
Using homoscedastic noise assumption:
$y^\star|\mathbf{x}^\star,\mu,k\sim\mathcal{N}(\mu(\mathbf{x}^\star),k(\mathbf{x}^\star,\mathbf{x}^\star)+\sigma_n^2)$
**New point**: Joint distribution of the observations $y_{1:n}$ and the noise-free prediction $f^\star$ at a test point $\mathbf{x}^\star$ as $\begin{bmatrix}\mathbf{y}\\f^\star\end{bmatrix}|\mathbf{x}^\star,\mathbf{x}_{1:n}\sim\mathcal{N}(\tilde{\mu},\tilde{\mathbf{K}})$ where $\tilde{\mu}\doteq\begin{bmatrix}\mu_A\\\mu(\mathbf{x}^\star)\end{bmatrix}$,
$\tilde{\mathbf{K}}\doteq\begin{bmatrix}\mathbf{K}_{AA}+\sigma_n^2\mathbf{I}&\mathbf{k}_{\mathbf{x}^\star,A}\\\mathbf{k}_{\mathbf{x}^\star,A}^\top&k(\mathbf{x}^\star,\mathbf{x}^\star)\end{bmatrix}$, $\mathbf{k}_{\mathbf{x},A}\doteq[k(\mathbf{x},\mathbf{x}_1)\dots k(\mathbf{x},\mathbf{x}_n)]^\top$

**GP posterior update**: $f|\mathbf{x}_{1:n},y_{1:n}\sim\mathcal{GP}(\mu',k')$
where $\mu'(\mathbf{x})\doteq\mu(\mathbf{x})+\mathbf{k}_{\mathbf{x},A}^\top(\mathbf{K}_{AA}+\sigma_n^2\mathbf{I})^{-1}(\mathbf{y}_A-\mu_A)$
and $k'(\mathbf{x},\mathbf{x}')\doteq k(\mathbf{x},\mathbf{x}')-\mathbf{k}_{\mathbf{x},A}^\top(\mathbf{K}_{AA}+\sigma_n^2\mathbf{I})^{-1}\mathbf{k}_{\mathbf{x}',A}$
- The posterior covariance can only decrease when conditioning on more data, and is independent of $\mathbf{y}_i$.
- GP posterior takes $\mathcal{O}(n^3)$ because of mat. inversion.

**Maximizing marginal likelihood**: optimizes $\theta$ across all realizations of $\mathbf{f}$ (somewhat reg., avoids overfitting):
$\hat{\theta}_{\text{MLE}}\doteq\arg\max_\theta p(y_{1:n}|\mathbf{x}_{1:n},\theta)$
$=\arg\max_\theta\int p(y_{1:n}|\mathbf{x}_{1:n},f,\theta)p(f|\theta)df$
- *Underfit* models: likelihood is mostly small as data cannot be well described; prior is large as there are "fewer" fns. to choose from.
- *Overfit* models: likelihood is large for "some" fns. but small for "most" fns.; prior is small, as probability mass has to be distributed among "more" fns.
Maximizing encourages trading between a large likelihood and large prior, as one product term will be small.

**For GP regression**: $y_{1:n}|\mathbf{x}_{1:n},\theta\sim\mathcal{N}(\mathbf{0},\mathbf{K}_{f,\theta}+\sigma_n^2\mathbf{I})$, write $\mathbf{K}_{\mathbf{y},\theta}\doteq\mathbf{K}_{f,\theta}+\sigma_n^2\mathbf{I}$, and obtain:
$\hat{\theta}_{\text{MLE}}=\arg\min_\theta\underbrace{1/2\cdot\mathbf{y}^\top\mathbf{K}_{\mathbf{y},\theta}^{-1}\mathbf{y}}_{\text{Goodness of fit}}+\underbrace{1/2\cdot\log\det(\mathbf{K}_{\mathbf{y},\theta})}_{\text{"Volume" of model class}}$
The loss can be expressed in closed-form with $\alpha\doteq\mathbf{K}^{-1}\mathbf{y}$:
$\frac{\partial}{\partial\theta_j}\log p(y_{1:n}|\mathbf{x}_{1:n},\theta)=\frac{1}{2}\text{tr}((\alpha\alpha^\top-\mathbf{K}_{\mathbf{y},\theta}^{-1})\frac{\partial\mathbf{K}_{\mathbf{y},\theta}}{\partial\theta})$
This optimization problem is, in general, non-convex.
- **GPs remain comp. tractable** even though they can model fns. over "infinite-dim." feat. spaces.
- For all $\mathbf{x}\in\mathcal{X}$ and $\mathbf{f}\in\mathcal{H}_k(\mathcal{X})$: $\mathbf{f}(\mathbf{x})=\langle f(\cdot),k(\mathbf{x},\cdot)\rangle_k$
Given kernel $k:\mathcal{X}\times\mathcal{X}\to\mathbb{R}$, its **RKHS** is space of fns.:
$\mathcal{H}_k(\mathcal{X})\doteq\{f(\cdot)=\sum_{i=1}^m\alpha_i k(\mathbf{x}_i,\cdot)|n\in\mathbb{N},\mathbf{x}_i\in\mathcal{X},\alpha_i\in\mathbb{R}\}$
Inner prod: $\langle f,g\rangle_k\doteq\sum_{i=1}^n\sum_{j=1}^{n'}\alpha_i\alpha_j'k(\mathbf{x}_i,\mathbf{x}_j')$
where $g(\cdot)=\sum_{j=1}^{n'}\alpha_j'k(\mathbf{x}'_j,\cdot)$. Induces norm
$\|\mathbf{f}\|_k=\sqrt{\langle\mathbf{f},\mathbf{f}\rangle_k}$ measuring smoothness/complexity.
- **Representer theorem:** Kernel $k$, $\lambda>0$, $f\in\mathcal{H}_k(\mathcal{X})$, and train data $\{(\mathbf{x}_i,f(\mathbf{x}_i))\}_{i=1}^n$. Let loss fn. $\mathcal{L}(f(\mathbf{x}_1),\dots,f(\mathbf{x}_n))\in\mathbb{R}\cup\{\infty\}$ depend on $f$ only through its eval. at train points. Then, any minimizer
  $\hat{f}\in\arg\min_{f\in\mathcal{H}_k(\mathcal{X})}\mathcal{L}(f(\mathbf{x}_1),\dots,f(\mathbf{x}_n))+\lambda\|f\|_k^2$ admits a repr. of form $\hat{f}(\mathbf{x})=\hat{\mathbf{a}}\mathbf{k}_{\mathbf{x},\{\mathbf{x}_i\}_{i=1}^n}=\sum_{i=1}^n\hat{\alpha}_i k(\mathbf{x},\mathbf{x}_i)$
- GP MAP is solution of reg. LR problem in RKHS of kernel fn.: $\hat{f}\doteq\arg\min_{f\in\mathcal{H}_k(\mathcal{X})}-\log(y_{1:n}|\mathbf{x}_{1:n},f)+\frac{1}{2}\|f\|_k^2$
- **GPs remain comp. tractable**

---

**Approximations**: GP need to invert mat ($\mathcal{O}(n^3)$).
- **Local method**: When sampling at $\mathbf{x}$ only condition on samples $\mathbf{x}'$, that are close: $|k(\mathbf{x},\mathbf{x}')|\geq\tau$ for some $\tau>0$. **Problem:** $\tau$ has to be chosen carefully: if $\tau$ is chosen too large, samples become essentially independent. Still expensive if "many" points are close.
- **Kernel approximation**: Construct a low dim. feat. map $\phi:\mathbb{R}^d\to\mathbb{R}^m$ s.t.: $k(\mathbf{x},\mathbf{x}')\approx\phi(\mathbf{x})^\top\phi(\mathbf{x}')$.
- Transforms function-space view (GP) back into a tractable weight-space view (BLR, $\mathcal{O}(nm^2+m^3)$).
- Can be done with **Random Fourier features**: a *stationary* kernel $k$ can be interpreted as fn. in one variable, and has an associated Fourier transform.
  **Bochner's Theorem** A continuous Kernel on $\mathbb{R}^d$ is p.s.d iff its Fourier transform $p(\omega)$ is non-negative.
  Rand. feat. map: $z_{\omega,b}(\mathbf{x})=\sqrt{2}\cos(\omega^\top\mathbf{x}+b)$, with $\omega^{(i)}\overset{\text{iid}}{\sim}p$ and $b^{(i)}\overset{\text{iid}}{\sim}\text{Unif}([0,2\pi])$. Inner product $z(x)^\top z(x')$ is unbiased estimator of $k(x-x')$. Error prob. decays exp. in dim. of Fourier feature space $m$.
- **Inducing point methods**: Idea is to summarize data around inducing pts. $U=\{\bar{\mathbf{x}}_1,\dots,\bar{\mathbf{x}}_n\}$. Let $\mathbf{f}\doteq[f(\mathbf{x}_1)\dots,f(\mathbf{x}_n)]^\top$, $f^\star\doteq[f(\mathbf{x}^\star)]$, $\mathbf{u}=[f(\hat{\mathbf{x}}_1)\dots f(\hat{\mathbf{x}}_n)]^\top$. Original GP recoverable with marginalization: $p(f^\star|\mathbf{f})=\int_{\mathbb{R}^k}p(f^\star,\mathbf{f}|\mathbf{u})p(\mathbf{u})d\mathbf{u}$. Approx. the joint prior, assuming $f^\star$, $\mathbf{f}$ are cond. indep. given $\mathbf{u}\sim\mathcal{N}(\mathbf{0},\mathbf{K}_{UU})$
  Train: $p(\mathbf{f}|\mathbf{u})\sim\mathcal{N}(\mathbf{f};\mathbf{K}_{AU}\mathbf{K}_{UU}^{-1}\mathbf{u},\mathbf{K}_{AA}-\mathbf{Q}_{AA})$
  Test: $p(f^\star|\mathbf{u})\sim\mathcal{N}(f^\star;\mathbf{K}_{\star U}\mathbf{K}_{UU}^{-1}\mathbf{u},\mathbf{K}_{\star\star}-\mathbf{Q}_{\star\star})$ w.
  $\mathbf{Q}_{ab}\doteq\mathbf{K}_{aU}\mathbf{K}_{UU}^{-1}\mathbf{K}_{Ub}$. $\mathbf{K}_{AA}$ represents the prior covar. and $\mathbf{Q}_{AA}$ represents covar. from inducing pts. Covar. mat. comp. is expensive; need to approx.:
  - **Subset of regressors (SoR)**: Forgets about all var. and covar. $q_{\text{SoR}}(\mathbf{f}|\mathbf{u})\doteq\mathcal{N}(\mathbf{f};\mathbf{K}_{AU}\mathbf{K}_{UU}^{-1}\mathbf{u},\mathbf{0})$
    $q_{\text{SoR}}(f^\star|\mathbf{u})\doteq\mathcal{N}(f^\star;\mathbf{K}_{\star U}\mathbf{K}_{UU}^{-1}\mathbf{u},\mathbf{0})$
  - **Fully independent training conditional (FITC)**: Keeps track of variances but forgets about covariance $q_{\text{FITC}}(\mathbf{f}|\mathbf{u})\doteq\mathcal{N}(\mathbf{f};\mathbf{K}_{AU}\mathbf{K}_{UU}^{-1}\mathbf{u},\text{diag}\{\mathbf{K}_{AA}-\mathbf{Q}_{AA}\})$
    $q_{\text{FITC}}(\mathbf{f}|\mathbf{u})\doteq\mathcal{N}(f^\star;\mathbf{K}_{\star U}\mathbf{K}_{UU}^{-1}\mathbf{u},\text{diag}\{\mathbf{K}_{\star\star}-\mathbf{Q}_{\star\star}\})$
  Comp. cost SoR/FITC is dom. by mat. inv. of $\mathbf{K}_{UU}$, so cubic in num. inducing pts. and linear in data pts.

# 4 Variational Inference

Idea: approximate true posterior distribution with a simpler posterior that is easy to sample:
$p(\theta|\mathbf{x}_{1:n},y_{1:n})=\frac{1}{Z}p(\theta,y_{1:n}|\mathbf{x}_{1:n})\approx q(\theta|\lambda)\doteq q_\lambda(\theta)$,
where $\lambda$ are params. of the **variational posterior** $q_\lambda$.

**Laplace approx.**: find a Gaussian approx. (i.e. second-order Taylor) of the posterior around its mode:
$q(\theta)\doteq\mathcal{N}(\theta;\hat{\theta},\Lambda^{-1})\propto\exp(\hat{\psi}(\theta))$, with $\hat{\theta}$ the mode (i.e. MAP estimate) and with $\mathbf{H}$ the Hessian:
$\Lambda\doteq-\mathbf{H}_\psi(\hat{\theta})=-\mathbf{H}_\theta\log p(\theta|\mathbf{x}_{1:n},y_{1:n})|_{\theta=\hat{\theta}}$.
Perform inference using the variations approximation:
$p(y^\star|\mathbf{x}^\star,\mathbf{x}_{1:n},y_{1:n})\approx\int p(y^\star|\mathbf{x}^\star,\theta)q_\lambda(\theta)d\theta$
$=\mathbb{E}_{\theta\sim q_\lambda}[p(y^\star|\mathbf{x}^\star,\theta)]$
- Matches shape of true posterior around its mode but may not represent it accurately elsewhere.
- Leads to extremely overconfident predictions, often unsuitable for approximate probabilistic inference.
- Preserves MAP point estimate as its mean.

**Jensen's inequality**: For convex fn. $g$ we have:
$g(\mathbb{E}[X])\leq\mathbb{E}[g(X)]$; if $h$ concave: $h(\mathbb{E}[X])\geq\mathbb{E}[h(X)]$

**Surprise** of an event with probability $u$: $S[u]\doteq-\log u$.
$S[u]$ is convex in $u$. For a discrete RV $X$: $S[p(x)]\geq0$.
Axiomatic characterization up to pos. const. factor:
- $S[u]>S[v]\Longrightarrow u<v$ (anti-monotonicity)
- $S$ continuous
- $S[uv]=S[u]+S[v]$ for independent events

The **cross-entropy** of $q$ relative to $p$ is:
$\text{H}[p\|q]\doteq\mathbb{E}_{x\sim p}[S[q(x)]]=\mathbb{E}_{x\sim p}[-\log q(x)]$.

---

**Entropy** of distr. $p$ is avg. surprise of samples from $p$:
$\text{H}[p]\doteq\mathbb{E}_{\mathbf{x}\sim p}[S[p(\mathbf{x})]]=\mathbb{E}_{\mathbf{x}\sim p}[-\log p(\mathbf{x})]$.
Can be negative if $p$ is discrete then $\text{H}[p]\geq0$.

**KL-div.:** measures additional expected surprise from observing samples from $p$ that is due to assuming (wrong) $q$ and which is not inherent in $p$ already.
$\text{KL}(p\|q)\doteq\text{H}[p\|q]-\text{H}[p]=\mathbb{E}_{\theta\sim p}[\log\frac{p(\theta)}{q(\theta)}]$
- $\text{KL}(p\|q)\geq0$; $\text{KL}(p\|q)=0$ iff $p=q$ almost surely
- There exist distr. $p$ and $q$ s.t. $\text{KL}(p\|q)\neq\text{KL}(q\|p)$
Note that: $\text{H}[p\|q]=\text{H}[p]+\text{KL}(p\|q)\geq\text{H}[p]$.
- $\text{KL}(\text{Bern}(p)\|\text{Bern}(q))=p\log\frac{p}{q}+(1-p)\log\frac{(1-p)}{(1-q)}$
- For $p\doteq\mathcal{N}(\mu_p,\Sigma_p)$ and $q\doteq\mathcal{N}(\mu_q,\Sigma_q)$: $\text{KL}(p\|q)=$
  $\frac{1}{2}(\text{tr}(\Sigma_q^{-1}\Sigma_p)+(\mu_p-\mu_q)^\top\Sigma_q^{-1}(\mu_p-\mu_q)-d+\log\frac{\det(\Sigma_q)}{\det(\Sigma_p)})$
- **Forward KL**: $q_1^\star\doteq\arg\min_{q\in\mathcal{Q}}\text{KL}(p\|q)$(mode avg., more conservative, yields more "desired" approx.)
- **Reverse KL**: $q_2^\star\doteq\arg\min_{q\in\mathcal{Q}}\text{KL}(q\|p)$ (greedily mode seeking, underestimate var., overconfident preds.)

**Evidence lower bound (ELBO)**, for given data $\mathcal{D}_n$:
$L(q,p;\mathcal{D}_n)\doteq\underbrace{\log p(y_{1:n}|\mathbf{x}_{1:n})}_{\text{const}}-\text{KL}(q\|p(\cdot|\mathbf{x}_{1:n},y_{1:n}))$
$=\underbrace{\mathbb{E}_{\theta\sim q}[\log p(y_{1:n}|\mathbf{x}_{1:n},\theta)]}_{\text{log-likelihood}}-\underbrace{\text{KL}(q\|p(\cdot))}_{\text{proximity to prior}}$
Max the ELBO coincides with min. reverse-KL. Since KL-div. is non-negative: $\log p(y_{1:n}|\mathbf{x}_{1:n})\geq L(q,p;\mathcal{D}_n)$
- Max. ELBO selects a var. distr. $q$ that is close to prior $p(\cdot)$ while also max. avg. data likelihood $p(y_{1:n}|\mathbf{x}_{1:n},\theta)$ for $\theta\sim q$. Contrast to MAP, which picks single mode $\theta$ that max. the likelihood and proximity to the prior.
- ELBO gradient is gen. **intractable** (use rep. trick).

**Reparam. trick:** Let $\epsilon\sim\phi$ be indep. of $\lambda$, $\mathbf{g}:\mathbb{R}^d\to\mathbb{R}^d$ be a diff. and inv. fn, $\theta=\mathbf{g}(\epsilon;\lambda)$, and $\mathbf{f}$ a *nice* fn. We get:
$q_\lambda(\theta)=\phi(\epsilon)\cdot|\det(\mathbf{D}_\epsilon\mathbf{g}(\epsilon;\lambda))|^{-1}$; $\mathbb{E}_{\theta\sim q_\lambda}[\mathbf{f}(\theta)]=\mathbb{E}_{\epsilon\sim\phi}[\mathbf{f}(\mathbf{g}(\epsilon;\lambda))]$.
- For ELBO $\nabla_\lambda\mathbb{E}_{\theta\sim q_\lambda}[\mathbf{f}(\theta)]=\mathbb{E}_{\epsilon\sim\phi}[\nabla_\lambda\mathbf{f}(\mathbf{g}(\epsilon;\lambda))]$.
- For Gaussian: $q_\lambda(\theta)\doteq\mathcal{N}(\theta;\mu,\Sigma)$; $\epsilon\sim\mathcal{N}(\mathbf{0},\mathbf{I})$, set:
  $\theta=\mathbf{g}(\epsilon;\lambda)\doteq\Sigma^{1/2}\epsilon+\mu$, then: $\phi(\epsilon)=q_\lambda(\theta)\cdot|\det(\Sigma^{1/2})|$
  and $\epsilon=\mathbf{g}^{-1}(\theta;\lambda)=\Sigma^{-1/2}(\theta-\mu)$

# 5 Bayesian Deep Learning

**Universal approx. theorem:** Any ANN with a single hidden layer (arbitrary width) and non-poly. activation fn. can approx. any cont. fn. to an arbitrary accuracy.
- **Hyperbolic tangent**: $\text{Tanh}(z)=\frac{\exp(z)-\exp(-z)}{\exp(z)+\exp(-z)}$
  $\nabla_z\text{Tanh}(z)=1-\text{Tanh}^2(z)$; $\text{Tanh}(z)=2\sigma(2z)-1$;
- **Softmax**: $\sigma_i(\mathbf{f})=\frac{\exp(f_i)}{\sum_{j=1}^c\exp(f_j)}$ (classification)
  $\nabla_z\sigma_i(\mathbf{f})=\sigma_i(\mathbf{f})(1-\sigma_i(\mathbf{f}))$    $\nabla_z\text{ReLU}(z)=\mathbf{1}_{\{z\geq0\}}$
- **Rectified linear unit**: $\text{ReLU}(z)\doteq\max\{z,0\}\in[0,\infty)$
For linear regression, min. MSE/CE **corresponds to** MLE under a Gaussian likelihood.
- **MSE loss**: $\ell_{\text{mse}}(\theta;\mathcal{D})\doteq\frac{1}{n}\sum_{i=1}^n(f(\mathbf{x}_i;\theta)-y_i)^2$
- **CE loss**: $\ell_{\text{ce}}(\theta;\mathcal{D})\doteq-\frac{1}{n}\sum_{i=1}^n\log q_\theta(\mathbf{y}_i|\mathbf{x}_i)$

**BNNs**: Gaussian prior on weights $\theta\sim\mathcal{N}(\mathbf{0},\sigma_p^2\mathbf{I})$, and Gaussian likelihood to describe how well data is described by the model: $y|\mathbf{x},\theta\sim\mathcal{N}(f(\mathbf{x};\theta),\sigma_n^2)$.
The **MAP estimate** is:
$\hat{\theta}_{\text{MAP}}=\arg\min_\theta\frac{1}{2\sigma_p^2}\|\theta\|_2^2+\frac{1}{2\sigma_n^2}\sum_{i=1}^n(y_i-f(\mathbf{x}_i;\theta))^2$.
Update rule: $\theta\leftarrow\theta(1-\frac{\eta_t}{\sigma_p^2})+\eta_t\frac{1}{\sigma_n^2}\sum_{i=1}^n\nabla\log p(y_i|\mathbf{x}_i,\theta)$.

**Also modeling heteroscedastic noise**: Use a neural network with 2 outputs $f_1,f_2$, and define:
$y|\mathbf{x},\theta\sim\mathcal{N}(\mu(\mathbf{x};\theta),\sigma^2(\mathbf{x};\theta))$ where $\mu(\mathbf{x};\theta)=f_1(\mathbf{x};\theta)$ and $\sigma^2(\mathbf{x};\theta)\doteq\exp(f_2(\mathbf{x};\theta))$. Likelihood term:
$\log p(y_i|\mathbf{x}_i,\theta)=\text{const}-\frac{1}{2}[\log\sigma^2(\mathbf{x}_i;\theta)+\frac{(y_i-\mu(\mathbf{x}_i;\theta))^2}{\sigma^2(\mathbf{x}_i;\theta)}]$.
- BNN learning and inference are **generally intractable** when the noise is not assumed to be homoscedastic and known. Thus, we need approx. inference.

- **Goal**: approx. true posterior $p(\theta|\mathcal{D})$ with simpler variational distr. $q_\lambda$ typically family of indep. Gaussians.
- Achieved by max. ELBO with SGD and rep. trick.
- We can approx. the predictive distr. by sampling from the variational posterior $p(y^*|\mathbf{x}^*,\mathbf{x}_{1:n},\mathbf{y}_{1:n})\approx$
$\mathbb{E}_{\theta\sim q_\lambda}[p(y^*|\mathbf{x}^*,\theta)]\approx\frac{1}{m}\sum_{i=1}^m p(y^*|\mathbf{x}^*,\theta^{(i)})$.
- VI in BNNs can be seen as avg. preds. of multiple NNs drawn acc. to the variational posterior $q_\lambda$.
- Using Monte Carlo samples estimate mean and var.:
$\mathbb{E}[y^*|\mathbf{x}^*,\mathbf{x}_{1:n},\mathbf{y}_{1:n}]\approx\frac{1}{m}\sum_{i=1}^m\mu(\mathbf{x}^*;\theta^{(i)})$
$\mathrm{Var}[y^*|\mathbf{x}^*,\mathbf{x}_{1:n},\mathbf{y}_{1:n}]\approx\mathbb{E}_\theta[\mathrm{Var}_{y^*}[y^*|\mathbf{x}^*,\theta]]+\mathrm{Var}_\theta[\mathbb{E}_{y^*}[y^*|\mathbf{x}^*,\theta]]$
$\approx\mathbb{E}_\theta[\sigma^2(\mathbf{x}^*;\theta)]+\mathrm{Var}_\theta[\mu(\mathbf{x}^*;\theta)]$
$\approx\underbrace{\frac{1}{m}\sum_{i=1}^m\sigma^2(\mathbf{x}^*;\theta^{(i)})}_{\text{aleatoric}}+\underbrace{\frac{1}{m-1}\sum_{i=1}^m(\mu(\mathbf{x}^*;\theta^{*(i)})-\bar\mu(\mathbf{x}^*))^2}_{\text{epistemic}}$

- **Dropout/Dropconnect** random select/omits vertices/edges of comp. graph. For valid interpretation of this as VI, also need to perform during inference.
- Dropout masks must overlap, making preds. highly corr., leading to underest. of epistemic uc. **Masksembles** mitigate by choosing fixed set of pre-defined masks.
- **Probabilistic ensembles:** learn $m$ different NN over random chosen subsets of train data for each network.
- **Evidence of val. set**; How well model desc. val. set?:
$\log p(y_{1:m}^{\text{val}}|\mathbf{x}_{1:m}^{\text{val}},\mathbf{x}_{1:n}^{\text{train}},y_{1:n}^{\text{train}})\approx\frac{1}{k}\sum_{j=1}^k\frac{1}{m}\sum_{i=1}^m\log p(y_i^{\text{val}}|\mathbf{x}_i^{\text{val}},\theta^{(j)})$
- **Frequency:** Proportion of samples in bin $m$ that belong to 1: $\mathrm{freq}(B_m)\doteq\frac{1}{|B_m|}\sum_{i\in B_m}\mathbf{1}\{\mathcal{Y}_i=1\}$
- **Confidence:** Avg. conf. of samples in bin $m$ belonging to 1: $\mathrm{conf}(B_m)\doteq\frac{1}{|B_m|}\sum_{i\in B_m}\mathbb{P}\{\mathcal{Y}_i=1|\mathbf{x}_i\}$

A model is **well-calibrated** if its confidence coincides with its acc. across many preds.: $\mathrm{freq}(B_m)\approx\mathrm{conf}(B_m)$
- **ECE:** $\ell_{\text{ECE}}\doteq\sum_{m=1}^M\frac{|B_m|}{n}|\mathrm{freq}(B_m)-\mathrm{conf}(B_m)|$
- **MCE:** Max instead of sum

## 6 Active Learning

**Cond. entropy:** $H[\mathbf{X}|\mathbf{Y}]\doteq\mathbb{E}_{\mathbf{y}\sim p(\mathbf{y})}[H[\mathbf{X}|\mathbf{Y}=\mathbf{y}]]$
$=\mathbb{E}_{(\mathbf{x},\mathbf{y})\sim p(\mathbf{x},\mathbf{y})}[-\log p(\mathbf{x}|\mathbf{y})]$

**Joint entropy:** $H[\mathbf{X},\mathbf{Y}]\doteq\mathbb{E}_{(\mathbf{x},\mathbf{y})\sim p(\mathbf{x},\mathbf{y})}[-\log p(\mathbf{x},\mathbf{y})]$
- $H[\mathbf{X}|\mathbf{Y}]\neq H[\mathbf{Y}|\mathbf{X}]$ in general; but $H[\mathbf{X},\mathbf{Y}]=H[\mathbf{Y},\mathbf{X}]$
- $H[\mathbf{X},\mathbf{Y}]=H[\mathbf{Y}]+H[\mathbf{X}|\mathbf{Y}]=H[\mathbf{X}]+H[\mathbf{Y}|\mathbf{X}]$
- $H[\mathbf{X}]=H[\mathbf{Y}|\mathbf{X}]+H[\mathbf{X}]-H[\mathbf{Y}]$ (Bayes Rule)
- $H[\mathbf{X}|\mathbf{Y}]\leq H[\mathbf{X}]$ (Gibbs; Information never hurts)
- $\Leftrightarrow 0\leq H[\mathbf{X}]-H[\mathbf{X}|\mathbf{Y}]=I(\mathbf{X};\mathbf{Y})$

**Mutual info:** $I(\mathbf{X};\mathbf{Y})\doteq H[\mathbf{X}]+H[\mathbf{Y}]-H[\mathbf{X},\mathbf{Y}]$
- $I(\mathbf{X};\mathbf{Y})=I(\mathbf{Y};\mathbf{X})=\mathbb{E}_{\mathbf{y}\sim p}[KL(p(\mathbf{x}|\mathbf{y})\|p(\mathbf{x}))]$
**Cond. mutual info:**
- $I(\mathbf{X};\mathbf{Y}\mid\mathbf{Z})=H[\mathbf{X}|\mathbf{Z}]-H[\mathbf{X}|\mathbf{Y},\mathbf{Z}]$
$=H[\mathbf{X},\mathbf{Z}]+H[\mathbf{Y},\mathbf{Z}]-H[\mathbf{Z}]-H[\mathbf{X},\mathbf{Y},\mathbf{Z}]$
$=I(\mathbf{X};\mathbf{Y},\mathbf{Z})-I(\mathbf{X};\mathbf{Z})$
- $I(\mathbf{X};\mathbf{Y}|\mathbf{Z})=I(\mathbf{Y};\mathbf{X}|\mathbf{Z})$
- $I(\mathbf{X};\mathbf{Y};\mathbf{Z})\doteq I(\mathbf{X};\mathbf{Y})-I(\mathbf{X};\mathbf{Y}|\mathbf{Z})$, so the "information never hurts" principle does not hold for MI. Information about $Z$ may reduce the MI between $\mathbf{X}$ and $\mathbf{Y}$.
- Given (discrete) fn. $F:\mathcal{P}(\mathcal{X})\to\mathbb{R}$, the **marginal gain** of $\mathbf{x}\in\mathcal{X}$ given $A\subseteq\mathcal{X}$ is: $\Delta_F(\mathbf{x}|A)\doteq F(A\cup\{\mathbf{x}\})-F(A)$.
- The fn. is **submodular** iff for any $\mathbf{x}\in\mathcal{X}$ and any $A\subseteq B\subseteq\mathcal{X}$: $F(A\cup\{\mathbf{x}\})-F(A)\geq F(B\cup\{\mathbf{x}\})-F(B)$ or equally $\Delta_F(\mathbf{x}|A)\geq\Delta_F(\mathbf{x}|B)$. Submodularity can be interpreted as notion of "concavity" for discrete fns.
- It is called **monotone** if $F(A)\leq F(B)$.

**Maximization objective:** monotone submodular function: $I(S)\doteq I(\mathbf{f}_S;\mathbf{y}_S)=H[\mathbf{f}_S]-H[\mathbf{f}_S\mid\mathbf{y}_S]$. $H[\mathbf{f}_S]$: uc in $\mathbf{f}_S$ before observing $\mathbf{y}_S$. $H[\mathbf{f}_S|\mathbf{y}_S]$ uc in $\mathbf{f}_S$ after observing $\mathbf{y}_S$. Max. MI is in general NP-hard.
- **Greedy:** Pick $\mathbf{x}_1$ through $\mathbf{x}_n$ individually by greedily finding the location with the maximal MI, this provides a $(1-1/e)$-approximation of the optimum.

- **Uncertainty sampling**: Have already picked $S_t=\{\mathbf{x}_1,\ldots,\mathbf{x}_t\}$; Solve the following:
$\mathbf{x}_{t+1}\doteq\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}\Delta_I(\mathbf{x}|S_t)=\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}I(f_\mathbf{x};y_\mathbf{x}|\mathbf{y}_{S_t})$.
Doesn't work with heteroscedastic noise: large aleatoric uc may dominate epistemic uc. In classification corresponds to selecting label that max. entropy of predicted label: $\mathbf{x}_{t+1}\doteq\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}H[y_\mathbf{x}|\mathbf{x}_{1:t},y_{1:t}]$.

**Bayesian active learning by disagreement (BALD):**
Identifies points $\mathbf{x}$ where models *disagree* about label $y_\mathbf{x}$ (each model is *confident* but predict different labels):
$\mathbf{x}_{t+1}\doteq\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}I(\theta;y_\mathbf{x}|\mathbf{x}_{1:t},y_{1:t})=$
$\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}H[y_\mathbf{x}|\mathbf{x}_{1:t},y_{1:t}]-\mathbb{E}_\theta[\mathbf{x}_{1:t},y_{1:t}H[y_\mathbf{x}|\theta]]$
- **Inductive learning** extract general rules from data. Typically, we can directly observe $f(\mathbf{x})$ at any $\mathbf{x}$.
- **Transductive learning** make best pred. at particular $\mathbf{x}^*$. Typically, cannot directly observe $f(\mathbf{x}^*)$. Require gen. $f(\mathbf{s})$ from the behavior of $f$ at other locations.

## 7 Bayesian Optimization

**Cumulative regret** associated with choices $\{\mathbf{x}_t\}_{t=1}^T$ is:
$R_T\doteq\sum_{t=1}^T\underbrace{(\max_\mathbf{x}f^*(\mathbf{x})-f^*(\mathbf{x}_t))}_{\text{instantaneous regret}}$.

Goal: Achieve **sublinear regret**: $\lim_{T\to\infty}R_T/T=0$ (requires balancing exploration and exploitation).

**Algorithm 9.2:** Bayesian optimization (with GPs)

initialize $f\sim\mathcal{GP}(\mu_0,k_0)$
for $t=1$ to $T$ do
  choose $\mathbf{x}_t=\arg\max_{\mathbf{x}\in\mathcal{X}}F(\mathbf{x};\mu_{t-1},k_{t-1})$
  observe $y_t=f(\mathbf{x}_t)+\epsilon_t$
  perform a probabilistic update to obtain $\mu_t$ and $k_t$

- Common to use an **acquisition fn.** to greedily pick the next point to sample based on the current model.
- **Upper confidence bound (UCB):**
$\mathbf{x}_{t+1}\doteq\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}\mu_t(\mathbf{x})+\beta_{t+1}\sigma_t(\mathbf{x})$, where
$\sigma_t(\mathbf{x})\doteq\sqrt{k_t(\mathbf{x},\mathbf{x})}$. If $\beta_t=0$ then UCB is purely exploitative; if $\beta_t\to\infty$, UCB recovers uc sampling. UCB fn. generally non-convex.

When choosing $\beta_t$ appropriately: $R_T=\mathcal{O}(\sqrt{T\gamma_T})$, with $\gamma_T\doteq\max_{S\subseteq\mathcal{X}}I(\mathbf{f}_S;\mathbf{y}_S)=\max_{S\subseteq\mathcal{X},|S|=T}\frac{1}{2}\log\det(\mathbf{I}+\sigma_n^{-2}\mathbf{K}_{SS})$
is the maximum information gain after $T$ rounds.
- Linear: $\gamma_T=\mathcal{O}(d\log T)$
- Gaussian: $\gamma_T=\mathcal{O}((\log T)^{d+1})$
- Matérn for $\nu>\frac{1}{2}$: $\gamma_T=\mathcal{O}(T^{\frac{d}{2\nu+d}}(\log T)^{\frac{2\nu}{2\nu+d}})$

**Thompson Sampling:** At time $t+1$, we sample a fn. $\tilde{f}_{t+1}\sim p(\cdot|\mathbf{x}_{1:t},y_{1:t})$ from our posterior distr. Then, we simply max. $\tilde{f}_{t+1}$, $\mathbf{x}_{t+1}\doteq\mathrm{argmax}_{\mathbf{x}\in\mathcal{X}}\tilde{f}_{t+1}(\mathbf{x})$.

## 8 Diffusion generative models

Let $\beta_t\in(0,1]$, $\bar\alpha_t=\prod_{s=1}^t\alpha_s$, and $\alpha_s=1-\beta_s$. Typically, $\beta_t$ is monotonically increases, which implies that $\bar\alpha_t\to 0$ and thus $x_T\to\mathcal{N}(0,\mathbf{I})$ for $T\to\infty$.
1. **Forward process:** Transform data points into (Gaussian) noise by using a fixed noising MC $q$:
$q(x_{1:T}|x_0)=\prod_{t=1}^T q(x_t|x_{t-1})$
$q(x_t|x_{t-1})=\mathcal{N}(x_t;\sqrt{1-\beta_t}x_{t-1},\beta_t\mathbf{I})$
$q(x_t|x_0)=\mathcal{N}(x_t;\sqrt{\bar\alpha_t}x_{t-1},(\bar\alpha-1)\mathbf{I})$
2. **Backward process:** Learn a denoising MC $p$ matching the reversed forward process.
$p_\lambda(x_{t-1}|x_t)=\mathcal{N}(x_{t-1};\mu_\lambda(x_t,t),\Sigma_\lambda(x_t,t))$
$p_\lambda(x_{0:T})=p_\lambda(x_T)\prod_{t=1}^T p_\lambda(x_{t-1}|x_t)$
$p_\lambda(x_0)=\int p_\lambda(x_{0:T})dx_{1:T}$ where $x_{1:T}$ latent vars.
3. **Generation:** Now generate novel data points by simulating the learned denoising MC $p$. Sample seq.:
(1) $x_1\sim p(X_1)$, (2) $x_2\sim p(X_2|X_1=x_1)$, ...

Note $p_\lambda(x_0)$ is intractable. Idea: use VI. ELBO:
$\log p_\lambda(x_0)\geq\log p_\lambda(x_0)-\mathcal{D}_{KL}(q(\cdot|x_0)\|p_\lambda(\cdot|x_0))$
$=\mathbb{E}_r[\log p_\lambda(x_T)-\sum_{t=2}^T\frac{1}{T}\log\frac{q(x_t|x_{t-1})}{p_\lambda(x_{t-1}|x_t)}-\log\frac{q(x_1|x_0)}{p_\lambda(x_0|x_1)}]$
$=\mathrm{const}+\mathbb{E}_r[-\sum_{t=2}^T\underbrace{\mathcal{D}_{KL}(q(\cdot|x_t,x_0)\|p_\lambda(\cdot|x_t))}_{L_t}+\underbrace{\log p_\lambda(x_0|x_1)}_{L_1}]$
with "$x_{1:T}\sim q(\cdot|x_0)$. Now optimize this via **stochastic VI** using closed-form expression of this loss/the KL-divergence term with const. var. schedule:
$D_{KL}(q(\cdot|x_t,x_0)\|p_\lambda(\cdot|x_t))=\frac{1}{2\sigma_t^2}\|\mu_t'(x_t,x_0)-\mu_\lambda(x_t,t)\|_2^2+\mathrm{const}$
with $\mu_t'(x_t,x_0)=\frac{\sqrt{\bar\alpha_t}\beta_t}{1-\bar\alpha_t}x_0+\frac{\sqrt{\alpha_t}(1-\bar\alpha_{t-1})}{(1-\bar\alpha_t)}x_t$

## 9 Markov Decision Processes (MDPs)

- A policy induces a MC $(X_t^\pi)_{t\in\mathbb{N}_0}$: $p^\pi(x'|x)\doteq$
$\mathbb{P}(X_{t+1}^\pi=x'|X_t^\pi=x)=\sum_{a\in A}\pi(a|x)p(x'|x,a)$.
- The **discounted payoff** from time $t$ is:
$G_t\doteq\sum_{m=0}^\infty\gamma^m R_{t+m}$, for $\gamma\in[0,1)$.
- **State value fn.:** $v_t^\pi\doteq\mathbb{E}_\pi[G_t|X_t=x,A_t=a]$ measures avg. discounted payoff from time $t$ starting from $x\in X$.
- **State-action value fn. (Q-fn.):** $q_t^\pi(x;a)\doteq$
$\mathbb{E}_\pi[G_t|X_t=x,A_t=a]=r(x,a)+\gamma\sum_{x'\in X}p(x'|x,a)\cdot v_{t+1}^\pi(x')$
measures avg. discounted payoff from time $t$ starting from $x\in X$ and with playing action $a\in A$.

**Bellman Expectation Equation:**
- $v^\pi(x)=r(x,\pi(x))+\gamma\mathbb{E}_{x'|x,\pi(x)}[v^\pi(x')]$
- If stochastic policy: $v^\pi(x)=\mathbb{E}_{a\sim\pi(x)}[q^\pi(x,a)]$
$q^\pi(x,a)=r(x,a)+\gamma\mathbb{E}_{x'|x,a}\mathbb{E}_{a'\sim\pi(x')}[q^\pi(x',a')]$
- For deterministic: $v^\pi(x)=q^\pi(x,\pi(x))$.
Can be used to find $v^\pi$ given policy $\pi$, by solving linear system of eq. in cubic time in size of state space. Can also be solved using fixed pt. iter.: $\mathbf{B}^\pi\mathbf{v}\doteq\mathbf{r}^\pi+\gamma\mathbf{P}^\pi\mathbf{v}$.
$\|v_t^\pi-v^\pi\|_\infty=\|\mathbf{B}^\pi v_{t-1}^\pi-\mathbf{B}^\pi v^\pi\|_\infty=\gamma\|v_{t-1}^\pi-v^\pi\|_\infty\leq\gamma\|v_0^\pi-v^\pi\|_\infty$

**Bellman's theorem:** A policy $\pi^*$ is optimal iff it's greedy w.r.t its own value fn. In other words, $\pi^*$ is optimal iff $\pi^*(x)$ is a distr. over set $\mathrm{argmax}_{a\in A}q^*(x,a)$.
- If for every state there is a unique action that max. the q-fn., $\pi^*$ is deter. and unique. $\pi^*(x)=\mathrm{argmax}_{a\in A}q^*(x,a)$.
- For finite MDPs, PI converges to $\pi^*$ in poly. num. of iter. Each step takes cubic comp. in the num. of states.

**Algorithm 10.14:** Policy iteration
initialize $\pi$ (arbitrarily)
repeat
  compute $v^\pi$
  compute $\pi_{v^\pi}$
  $\pi\leftarrow\pi_{v^\pi}$
until converged

Monotonic improvement of PI:
- $v^{\pi^{t+1}}(x)\geq v^{\pi^t}(x)$ for all $x\in X$
- $v^{\pi^{t+1}}(x)>v^{\pi^t}(x)$ for at least one $x\in X$, unless $v^{\pi^t}\equiv v^*$

**Algorithm 10.17:** Value iteration
initialize $v(x)\leftarrow\max_{a\in A}r(x,a)$ for each $x\in X$
for $t=1$ to $\infty$ do
  $v(x)\leftarrow(\mathbf{B}^*v)(x)=\max_{a\in A}q(x,a)$ for each $x\in X$
choose $\pi_v$

- VI converges to an optimal policy, as $v^*$ and $q^*$ are a fixed-points of the Bellman update $\mathbf{B}^*$.
- For any $\epsilon>0$, VI converges to an $\epsilon$-optimal solution in poly time. However, unlike PI, VI does not generally reach the exact optimum in a finite num. of iter.

POMDP: Markov process, with **observations** $Y$, and **observation probs.** $o(y|x)\doteq\mathbb{P}(Y_t=y|X_t=x)$. Hard to solve in gen., can conv. to MDP with larger state space.

## 10 Tabular Reinforcement Learning

Markovian property of the underlying MDP:
$X_{t+1}\perp X_{t'+1}'|X_t,X_{t'+1}',A_t,A_{t'}\quad R_t\perp X_{t'}'|R_t,X_{t'+1}',A_t,A_{t'}$

**Bootstrapping:** approx. a true quantity by using an empirical quantity, which itself is constructed using samples from the true quantity that is to be approx.

**Robbins–Montro (RM) conditions:** for a sequence $(\alpha_t)_{t\in\mathbb{N}_0}$ if: $\alpha_t\geq 0$, $\sum_{t=0}^\infty\alpha_t=\infty$, $\sum_{t=0}^\infty\alpha_t^2<\infty$.

**For model-based approaches MLE yields:**
$\hat{p}(x'|x,a)=\frac{N(x'|x,a)}{N(a|x)}$ and $\hat{r}(x,a)=\frac{1}{N(a|x)}\sum_{t=0,x_t=x,a_t=a}^n r_t$
Both unbiased as they correspond to a sample mean.
- $N(x'|x,a)$ num. trans. from $x$ to $x'$ when play $a$.
- $N(a|x)$ nums trans. from $x$ and play $a$.

**Greedy in the limit with inf. exploration (GLIE):**
1. All state-action pairs are explored infinitely many times: $\lim_{t\to\infty}N_t(x,a)=\infty$
2. The policy converges to a greedy policy:
$\lim_{t\to\infty}\pi_t(a|x)=\mathbf{1}\{a=\mathrm{argmax}_{a'\in A}Q_t^*(x,a')\}$

- $\epsilon$-greedy: Ignores all past experience. Will eventually converge. Will eventually converge. GLIE with prob. 1 if $(\epsilon_t)_{t\in\mathbb{N}_0}$ satisfies the RM conditions (e.g., $\epsilon_t=1/t$).

**Softmax/Boltzmann exploration:** alt. to $\epsilon$-greedy
$\pi_\lambda(a|x)\propto\exp(\frac{1}{\lambda}Q^*(x,a))$ (Gibbs). For $\lambda\to 0$ greedily max. Q-fn. For $\lambda\to\infty$ uniform rand. exploration.

**Algorithm 11.6:** $R_{\max}$ algorithm

add the fairy-tale state $x^*$ to the Markov process
set $\hat{r}(x,a)=R_{\max}$ for all $x\in X$ and $a\in A$
set $\hat{p}(x^*|x,a)=1$ for all $x\in X$ and $a\in A$    *On-policy, Model-based*
compute the optimal policy $\hat\pi$ for $\hat{p}$ and $\hat{p}$
for $t=0$ to $\infty$ do
  execute policy $\hat\pi$ (for some number of steps)
  for each visited state-action pair $(x,a)$, update $\hat{r}(x,a)$
  estimate transition probabilities $\hat{p}(x'\mid x,a)$
  after observing "enough" transitions and rewards, recompute the
    optimal policy $\hat\pi$ according the current model $\hat{p}$ and $\hat{r}$.

- Optimism in the face of uc. Init. with max reward.
- Every $T$ steps, with high prob., either obtains near-optimal rew.; or visits one unknown state-action pair.
- With prob. at least $1-\delta$, $R_{\max}$ reaches $\epsilon$-optimal $\pi$ in poly. num. steps $|X|$, $|A|$, $T$, $1/\epsilon$, $1/\delta$, and $R_{\max}$.

**Algorithm 11.9:** Temporal-difference (TD) learning
initialize $V^\pi$ arbitrarily (e.g., as $\mathbf{0}$)    *On-policy, Model-free*
for $t=0$ to $\infty$ do
  follow policy $\pi$ to obtain the transition $(x,a,r,x')$
  $V^\pi(x)\leftarrow(1-\alpha_t)V^\pi(x)+\alpha_t(r+\gamma V^\pi(x'))$

- If $\alpha_t$ satisfies RM conditions and all state-action pairs are chosen inf. often, then $V^\pi$ conv. to $v^\pi$ w. prob 1.
- For estimates $V^\pi$ to converge true $v^\pi$, the transitions that are used for the estimation must follow policy $\pi$.

**SARSA:** Same as TD but estimate $Q$ with update:
$Q^\pi(x,a)\leftarrow(1-\alpha_t)Q^\pi(x,a)+\alpha_t(r+\gamma Q^\pi(x',a'))$
Same convergence guarantees as TD.    *On-policy, Model-free*

**Algorithm 11.12:** Q-learning
initialize $Q^*(x,a)$ arbitrarily (e.g., as $\mathbf{0}$)
for $t=0$ to $\infty$ do    *Off-policy, Model-free*
  observe the transition $(x,a,r,x')$
  $Q^*(x,a)\leftarrow(1-\alpha_t)Q^*(x,a)+\alpha_t(r+\gamma\max_{a'\in A}Q^*(x',a'))$

- If $\alpha_t$ satisfies RM cond. and all state-action pair are visited inf. often, then $Q^*$ conv. to $q^*$ with prob. 1.
- With prob. at least $1-\delta$, conv. to $\epsilon$-optimal policy in num. steps poly. in $\log|X|$, $\log|A|$, $\frac{1}{\epsilon}$ and $\log\frac{1}{\delta}$.

**Optimistic Q-learning:** Similar to $R_{\max}$ Init.
$Q^*(x,a)=V_{\max}\prod_{t=1}^{T_{\text{init}}}(1-\alpha_t)^{-1}$ w.
$V_{\max}=\frac{R_{\max}}{1-\gamma}\geq\max_a q^*(x,a)$. With prob. at least $1-\delta$, $\epsilon$-optimal $\pi$ after num. steps poly. in $|X|$, $|A|$, $\frac{1}{\epsilon}$, $\log\frac{1}{\delta}$, and $R_{\max}$ where $T_{\text{init}}$ is upper bounded by a poly. in same coeff. If $T_{\text{init}}$ large enough, conv. quickly to $\pi^*$.

## 11 Model-free Reinforcement Learning

In tab. methods: Storing val. fn., need at least $O(|\mathcal{X}|)$ space. Storing Q-fn, even need $O(|\mathcal{X}|\cdot|\mathcal{A}|)$ space. Time req. to compute value fn. for every state-action pair exactly grows poly. in size of state-action space. Can view TD-/Q-learning as SGD on the squared loss: $\ell(\theta;x,r,x')\doteq\frac{1}{2}(r+\gamma\theta^{\text{old}}(x')-\theta(x))^2$ and learn param. approx. of $V(\mathbf{x};\theta)$ or $Q(\mathbf{x},\mathbf{a};\theta)$ using Monte Carlo est. and bootstrapping.

Model-free (TD-/Q-learn.) are usually **sample ineff.**:
- Bootstrapping leads to "(initially) incorrect" and "unstable" targets of the optimization.
- Monte Carlo est. with single sample leads to large var.
**Q-learning with fn. approx.:** (1) Observe $\mathbf{x}'$ and $r$ from picking $a$ in $\mathbf{x}$. (2) Update $\theta\leftarrow\theta+\alpha_t\delta_B\phi(\mathbf{x},\mathbf{a})$, where $\delta_B\doteq r+\gamma\max_{\mathbf{a}'\in A}Q^*(\mathbf{x}',\mathbf{a}';\theta^{\text{old}})-Q^*(\mathbf{x},\mathbf{a};\theta)$.
- In the tabular setting, this is identical to Q-learning.
- Converges to the true Q-function $q^*$.
"Tricks of the trade" to improve SGD:
- **Stabilizing opti. targets:** Bootstrapping est. changes after each iteration, leading to stability issues. **DQN** updates NN used for bootstrapping infrequently and maintains const. opti. target across multiple episodes. E.g. clone: hanging/online NN and fixed/target NN.
- **Max. bias:** Estimates $Q^*$ are noisy (biased) estimates of $q^*$. **DDQN**: instead of picking optimal action w.r.t. old network, it picks w.r.t. new network:

**Policy val. fn.:** measures discounted payoff of policy:
$J(\pi)\doteq\mathbb{E}_\pi[G_0]=\mathbb{E}_\pi[\sum_{t=0}^\infty\gamma^t R_t]$, and bounded ver.:
$J_T(\pi)\doteq\mathbb{E}_\pi[G_{0:T}]=\mathbb{E}_\pi[\sum_{t=0}^{T-1}\gamma^t R_t]$.    Non-convex
**Score grad. est.:** $\nabla_\varphi\mathbb{E}_{\tau\sim\Pi_\varphi}[G_{0:T}]=\mathbb{E}[G_{0:T}\nabla_\varphi\log\Pi_\varphi(\tau)]$
Typically var. of est. very large. Reduce w. **baselines**:
$\mathbb{E}[G_{0:T}\nabla_\varphi\log\Pi_\varphi(\tau)]=\mathbb{E}[(G_{0:T}-b)\nabla_\varphi\log\Pi_\varphi(\tau)]$

**Algorithm 12.8:** REINFORCE algorithm
initialize policy weights $\varphi$    *On-policy, Model-free*
repeat
  generate an episode (i.e., rollout) to obtain trajectory $\tau$
  for $t=0$ to $T-1$ do
    set $g_{t:T}$ to the downstream return from time $t$
    $\varphi\leftarrow\varphi+\eta\gamma^t g_{t:T}\nabla_\varphi\log\pi_\varphi(a_t\mid x_t)$
until converged

- SGD with score grad. est. and downstream returns.
- Not guaranteed to find an optimal policy. Can get stuck in local optima even for very small domains.

**Advantage fn.:** $a^\pi(\mathbf{x},\mathbf{a})\doteq q^\pi(\mathbf{x},\mathbf{a})-v^\pi(\mathbf{x})$
$=q^\pi(\mathbf{x},\mathbf{a})-\mathbb{E}_{\mathbf{a}'\sim\pi(\mathbf{x})}[q^\pi(\mathbf{x},\mathbf{a}')]$
- $\pi$ is optimal $\Longleftrightarrow\forall\mathbf{x}\in\mathcal{X},\mathbf{a}\in\mathcal{A}:a^\pi(\mathbf{x},\mathbf{a})\leq 0$

**Policy gradient theorem:** Max. $J(\varphi)$ corresponds to incr. the prob. of actions with large and decr. the prob. of actions with small value, taking into account how often the resulting policy visits certain states.
$\nabla_\varphi J(\varphi)=\sum_{t=0}^\infty\mathbb{E}_{\mathbf{x}_t,\mathbf{a}_t}[\gamma^t q_\varphi^\pi(\mathbf{x}_t,\mathbf{a}_t)\nabla_\varphi\pi_\varphi(\mathbf{a}_t,\mathbf{x}_t)]$

**Algorithm 12.11:** Online actor-critic
initialize parameters $\varphi$ and $\theta$    *On-policy, Online, Model-free*
repeat
  use $\pi_\varphi$ to obtain transition $(x,a,r,x')$ and the next $a$
    $a'\sim\pi_\varphi(\cdot\mid x')$
  $\delta=r+\gamma Q(x',a';\theta)-Q(x,a;\theta)$
  // actor update
  $\varphi\leftarrow\varphi+\eta Q(x,a;\theta)\nabla_\varphi\log\pi_\varphi(a\mid x)$
  // critic update
  $\theta\leftarrow\theta+\eta\delta\nabla_\theta Q(x,a;\theta)$
until converged

- Use SARSA for learning critic; SGD for gradient est.
- Actor is not guaranteed to improve

- **TRPO:** KL-constraint forces monotonic improvement via trust regions (on-pol., m-free, pol-gradient) ● **PPO:** Clipped obj. stabilizes training by limiting update size (on-pol., m-free, AC) ● **GRPO:** Group-relative advantages remove need for separate critic (on-pol., m-free, critic-less) ● **DDPG:** Det. grads. enable cont. control with target networks (off-pol., m-free, AC) ● **SAC:** Max. entropy obj. balances exploration and performance (off-pol., m-free) ● **DPO:** Bypasses explicit reward modeling (offline, m-free, ref-based) ● **PETS:** Ensembles plus MPC handle uncertainty in planning. (m-based)
- **UCRL:** selects optimistic MDPs for efficient exploration (m-based, exploration, o.i.t.f.o.uc) ● **H-UCRL:** hallucinated transitions ensure safe, robust exploration. (m-based, opt. o.i.t.f.o.uc, hierarchical)