

BI-ML2.21 přednáška 1

Daniel Vašata

FIT ČVUT

21. 2. 2024

Autor: Daniel Vašata.

Problémy, návrhy apod. hlaste v [GitLabu](#).

Verze souboru: 21. února 2024 12:43.

Podmínky získání zápočtu

- **Zápočet bude postaven na vypracovávání domácích úkolů.**
- Úkoly budou během semestru **dva**, každý za **max. 25 bodů**.
- Domácí úkoly budete vypracovávat v jazyce Python ve formátu Jupyter notebook (.ipynb).
- Zadání a podrobné instrukce k vypracování a odevzdání najdete na stránkách předmětu:

 courses.fit.cvut.cz/BI-ML2/ 

- **K získání zápočtu je třeba získat alespoň 25 bodů z 50.**

Podmínky složení zkoušky

- Zkouška bude **ústní**.
- Každý student dostane **dvě otázky z předem zveřejněného seznamu**.
- Z každé otázky můžete získat **až 25 bodů**.
- Celkem tedy můžete ze zkoušky získat **až 50 bodů**. Není žádný minimální nutný počet bodů, který je třeba ze zkoušky získat.
- Pokud student/ka u zkoušky prokáže zásadní neznalost, může zkoušející použít **právo veta a zkoušku ukončit jako neúspěšnou**.
- V případě úspěchu u zkoušky se výsledná známka odvodí ze součtu bodů ze semestru a ze zkoušky.

O čem to všechno vlastně bude?

Budeme rozšiřovat znalosti získané v kurzu **BI-ML1**.

Zejména se budeme zabývat:

- Dalšími důležitými metodami supervizovaného učení
- Metodami redukce dimenzionality
- Neuronovými sítěmi
- Posilovaným učením

Doporučená literatura:

- Hastie T., Tibshirani R., Friedman J.: **The Elements of Statistical Learning**. Springer, 2009.
- Murphy K. P.: **Machine Learning, A Probabilistic Perspective**. MIT Press, 2012.
- Goodfellow I., Bengio Y., Courville A.: **Deep Learning**. MIT Press, 2016.
- Sutton R. S., Barto A. G.: **Reinforcement Learning**. MIT Press, 2018.

Co bude v dnešní přednášce

- Opakování modelu lineární a hřebenové regrese
- Opakování lineárního modelu bazových funkcí
- Duální reprezentace optimalizační úlohy pro trénování
- Diskuse jádrového triku a ukázky jádrových funkcí

Opakování lineární regrese

Začneme připomenutím lineární regrese.

Model lineární regrese

Hodnota vysvětlované proměnné Y v bodě o hodnotách x_1, \dots, x_p příznaků X_1, \dots, X_p je

$$Y = w_0 + w_1 x_1 + \dots + w_p x_p + \varepsilon,$$

kde w_0, \dots, w_p jsou neznámé parametry a ε je náhodná veličina pro kterou platí $E\varepsilon = 0$.

Poznámky:

- Zavedeme-li nový konstantní příznak $X_0 = x_0 = 1$ a vektorové značení

$$\mathbf{x} = (x_0, x_1, \dots, x_p)^T \quad \text{a} \quad \mathbf{w} = (w_0, w_1, \dots, w_p)^T,$$

můžeme zkráceně psát

$$Y = \mathbf{w}^T \mathbf{x} + \varepsilon.$$

- Naším cílem je **odhadnout** neznámé váhy \mathbf{w} pomocí $\hat{\mathbf{w}}$ a pak **predikovat** Y jako

$$\hat{Y} = \hat{\mathbf{w}}^T \mathbf{x}.$$

- Ze statistického pohledu je predikce \hat{Y} **bodovým odhadem** $EY = \mathbf{w}^T \mathbf{x}$.

Model pro trénovací množinu

Trénovací množina je tvořena N dvojicemi $(Y_1, \mathbf{x}_1), \dots, (Y_N, \mathbf{x}_N)$, které jsou **nezávislé** a pocházejí ze stejného rozdělení, tj.

$$Y_i = \mathbf{w}^T \mathbf{x}_i + \varepsilon_i,$$

kde $\varepsilon_1, \dots, \varepsilon_N$ jsou nezávislé a $E \varepsilon_i = 0$.

Toto zapisujeme maticově jako

$$\mathbf{Y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon},$$

kde

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} 1 & x_{1;1} & \cdots & x_{1;p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N;1} & \cdots & x_{N;p} \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_N \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

s $E \boldsymbol{\varepsilon} = \mathbf{0}$.

Metoda nejmenších čtverců

- Při trénování minimalizujeme residuální součet čtverců

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (Y_i - \mathbf{x}_i^T \mathbf{w})^2 = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2.$$

- Minimum je určeno řešením **normální rovnice**

$$\mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{0},$$

která odpovídá podmínce $\nabla \text{RSS}(\mathbf{w}) = \mathbf{0}$.

- Za předpokladu, že je matice $\mathbf{X}^T \mathbf{X}$ regulární, existuje jediné řešení minimalizující $\text{RSS}(\mathbf{w})$,

$$\hat{\mathbf{w}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

- Predikce v bodě \mathbf{x} je potom $\hat{Y} = \mathbf{x}^T \hat{\mathbf{w}}_{\text{OLS}}$.
- Když $\mathbf{X}^T \mathbf{X}$ není regulární můžeme použít [Mooreovu-Penroseovu pseudoinverzní matici](#) $(\mathbf{X}^T \mathbf{X})^+$ místo $(\mathbf{X}^T \mathbf{X})^{-1}$ nebo přidat regularizační člen k $\text{RSS}(\mathbf{w})$ což vede na [hřebenovou regresi](#).

Hřebenová regrese

- Ve **hřebenové regresi** minimalizujeme **regularizovaný reziduální součet čtverců**

$$\text{RSS}_\lambda(\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \sum_{i=1}^p w_i^2,$$

který závisí na hyperparametru $\lambda \geq 0$.

- Zavedeme-li matici

$$\mathbf{I}' = \begin{pmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{p+1, p+1},$$

můžeme psát

$$\text{RSS}_\lambda(\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \mathbf{w}^T \mathbf{I}' \mathbf{w}.$$

- **Normální rovnice** je potom

$$\mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X} \mathbf{w} - \lambda \mathbf{I}' \mathbf{w} = \mathbf{0}.$$

- Protože $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}'$ je regulární pro každé $\lambda > 0$, existuje jednoznačné řešení

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}')^{-1} \mathbf{X}^T \mathbf{Y}.$$

- Hřebenová regrese pro $\lambda = 0$ dává obyčejnou metodu nejmenších čtverců, $\hat{\mathbf{w}}_0 = \hat{\mathbf{w}}$.

Lineární model báзовých funkcí

Doposud jsme byli schopni modelovat pouze lineární funkci ve vstupních proměnných.

Základní rozšíření spočívá v nahrazení původních příznaků jejich transformovanými variantami.

Označme jako \mathcal{X} prostor všech možných hodnot vektoru příznaků $\mathbf{X} = (X_1, \dots, X_p)^T$. Typicky, $\mathcal{X} = \mathbb{R}^p$ nebo alespoň $\mathcal{X} \subset \mathbb{R}^p$.

Pro $M \in \mathbb{N}$ uvažujme M **lineárně nezávislých** funkcí $\varphi_1, \dots, \varphi_M$ z \mathcal{X} do \mathbb{R} . Tyto funkce nazývané **báзовé funkce** (angl. **basis functions**) představují **transformace** původních příznaků X_1, \dots, X_p do **nového** M -rozměrného **příznakového prostoru**.

Jako model pro Y nyní použijme lineární model v tomto novém příznakovém prostoru.

Lineární model báзовých funkcí

Vysvětlovaná proměnná Y v bodě $\mathbf{x} = (x_1, \dots, x_p)^T$ hodnot vektoru příznaků \mathbf{X} je určena vztahem

$$Y = w_1\varphi_1(\mathbf{x}) + \dots w_M\varphi_M(\mathbf{x}) + \varepsilon = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + \varepsilon,$$

kde $\boldsymbol{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^M$ je vektorová funkce definována jako $\boldsymbol{\varphi}(\mathbf{x}) = (\varphi_1(\mathbf{x}), \dots, \varphi_M(\mathbf{x}))^T$ pro všechna $\mathbf{x} \in \mathcal{X}$, $\mathbf{w} = (w_1, \dots, w_M)^T$ je vektor neznámých parametrů, a ε je náhodná veličina s $E\varepsilon = 0$.

Pro jednoduchost při pozdějších úvahách nyní neuvažujeme intercept (který můžeme vytvořit vhodnou volbou jedné z báзовých funkcí).

Odhad v lineárním modelu bazových funkcí

- Model pro trénovací množinu tvořenou N dvojicemi $(Y_1, \mathbf{x}_1), \dots, (Y_N, \mathbf{x}_N)$ můžeme ve vektorovém tvaru zapsat jako

$$\mathbf{Y} = \Phi \mathbf{w} + \varepsilon,$$

kde

$$\Phi = \begin{pmatrix} \varphi(\mathbf{x}_1)^T \\ \vdots \\ \varphi(\mathbf{x}_N)^T \end{pmatrix} = \begin{pmatrix} \varphi_1(\mathbf{x}_1) & \cdots & \varphi_M(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \varphi_1(\mathbf{x}_N) & \cdots & \varphi_M(\mathbf{x}_N) \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_N \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{pmatrix}.$$

- Obecně budeme minimalizovat (pro $\lambda = 0$ máme metodu nejmenších čtverců)

$$\text{RSS}_\lambda(\mathbf{w}) = \|\mathbf{Y} - \Phi \mathbf{w}\|^2 + \lambda \mathbf{w}^T \mathbf{w}.$$

- Normální rovnice je

$$\Phi^T \mathbf{Y} - \Phi^T \Phi \mathbf{w} - \lambda \mathbf{w} = \mathbf{0}.$$

- Pro $\lambda > 0$ existuje jednoznačné řešení

$$\hat{\mathbf{w}}_\lambda = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{Y}.$$

- Predikce hodnoty Y v bodě \mathbf{x} je potom určena vztahem

$$\hat{Y} = \hat{\mathbf{w}}_\lambda^T \varphi(\mathbf{x}).$$

Bázové funkce

Mezi obvyklé volby bazových funkcí patří:

- $\varphi(\mathbf{x}) = x_i$ – přímo jednotlivé příznaky, odpovídají původnímu lineárnímu modelu.
- $\varphi(\mathbf{x}) = x_i^2, \varphi(\mathbf{x}) = x_k x_\ell$ – mocniny příznaků a jejich různé součiny, odpovídá polynomiální regresi.
- $\varphi(\mathbf{x}) = \log(x_i), \sqrt{x_i}, \sin(x_i)$ atd. – nelineární transformace jednotlivých příznaků.
- $\varphi(\mathbf{x}) = \mathbb{1}_{(a,b)}(x_i)$, where $\mathbb{1}_A(x) = 1$ if $x \in A$ a $\mathbb{1}_A(x) = 0$ if $x \notin A$ – indikátory množin. Umožňují rozdělení prostoru příznaků na kousky a následné modelování v každém kousku zvlášť.
- $\varphi(\mathbf{x}) = h(\|\mathbf{x} - \mathbf{x}_i\|)$, kde \mathbf{x}_i je i -tý trénovací bod a h je nějaká funkce – tzv. **radiální bazové funkce** centrované v bodech trénovací množiny.

Různé volby umožňují získat dostatečně flexibilní model pro Y .

Pokud nemáme žádné speciální znalosti o systému, který modelujeme, typicky na počátku volíme velké množství bazových funkcí a používáme hřebenovou regresi, případně jinou formu regularizace.

Duální reprezentace

Lineární modely bazových funkcí pro regresi a klasifikaci můžeme v mnoha případech přeformulovat do **duální reprezentace**, ve které se bazové funkce vyskytují pouze implicitně a jsou určeny pomocí tzv. jádrových funkcí.

Začneme s regresním lineárním modelem bazových funkcí, kdy minimalizujeme

$$\text{RSS}_\lambda(\mathbf{w}) = \|\mathbf{Y} - \Phi\mathbf{w}\|^2 + \lambda\mathbf{w}^T\mathbf{w}.$$

Uvažujme nyní hodnoty \mathbf{w} ve tvaru

$$\mathbf{w} = \Phi^T\boldsymbol{\alpha} \quad \text{kde} \quad \boldsymbol{\alpha} \in \mathbb{R}^N,$$

tj. omezení \mathbf{w} na podprostor \mathbb{R}^M generovaný vektory $\varphi(x_1), \dots, \varphi(x_N)$.

Dosazením do $\text{RSS}_\lambda(\mathbf{w})$ dostaneme

$$\text{RSS}_\lambda(\boldsymbol{\alpha}) = \|\mathbf{Y} - \Phi\Phi^T\boldsymbol{\alpha}\|^2 + \lambda\boldsymbol{\alpha}^T\Phi\Phi^T\boldsymbol{\alpha}.$$

Duální reprezentace

Jelikož je $\text{RSS}_\lambda(\alpha)$ určeno zúžením w na $w(\alpha) = \Phi^T \alpha$ dostáváme

$$\min_{\alpha} \text{RSS}_\lambda(\alpha) = \min_{\alpha} \text{RSS}_\lambda(w(\alpha)) \geq \min_w \text{RSS}_\lambda(w).$$

Následující věta ukazuje, že mezi hodnotami obou minim platí rovnost.

Věta

Pro $\lambda > 0$ platí

$$\min_{\alpha} \text{RSS}_\lambda(\alpha) = \min_w \text{RSS}_\lambda(w).$$

Navíc, pokud w^ minimalizuje $\text{RSS}_\lambda(w)$, tak $\alpha^* = \frac{1}{\lambda}(Y - \Phi w^*)$ minimalizuje $\text{RSS}_\lambda(\alpha)$.*

Na druhou stranu, pokud α^ minimalizuje $\text{RSS}_\lambda(\alpha)$, tak $w^* = \Phi^T \alpha^*$ minimalizuje $\text{RSS}_\lambda(w)$.*

Minimalizace $\text{RSS}_\lambda(\alpha)$ vzhledem k α je tedy ve skutečnosti **ekvivalentní** původní minimalizaci $\text{RSS}_\lambda(w)$ vzhledem w .

Duální reprezentace

Důkaz.

Je-li w^* argumentem minima $\text{RSS}_\lambda(w)$, pak splňuje normální rovnici odpovídající $\nabla \text{RSS}_\lambda(w^*) = 0$:

$$\Phi^T(Y - \Phi w^*) - \lambda w^* = 0. \quad (1)$$

Tudíž $w^* = \frac{1}{\lambda} \Phi^T(Y - \Phi w^*)$ a pro $\alpha^* = \frac{1}{\lambda}(Y - \Phi w^*)$ dostáváme $w^* = \Phi^T \alpha^*$. Z toho plyne $\text{RSS}_\lambda(\alpha^*) = \text{RSS}_\lambda(w^*)$ což znamená

$$\min_{\alpha} \text{RSS}_\lambda(\alpha) \leq \text{RSS}_\lambda(\alpha^*) = \text{RSS}_\lambda(w^*) = \min_w \text{RSS}_\lambda(w).$$

Protože $\min_{\alpha} \text{RSS}_\lambda(\alpha) \geq \min_w \text{RSS}_\lambda(w)$, ukázali jsme $\min_{\alpha} \text{RSS}_\lambda(\alpha) = \min_w \text{RSS}_\lambda(w)$ a tedy i první implikaci.

Pro obrácený směr uvažujme, že α^* minimalizuje $\text{RSS}_\lambda(\alpha)$. Potom také splní normální rovnici odpovídající $\nabla \text{RSS}_\lambda(\alpha^*) = 0$:

$$\Phi \Phi^T(Y - \Phi \Phi^T \alpha^*) - \lambda \Phi \Phi^T \alpha^* = 0.$$

Nyní stačí ukázat, že w^* sestrojené jako $w^* = \Phi^T \alpha^*$ splní rovnici (1).

$$\begin{aligned} \|\Phi^T(Y - \Phi w^*) - \lambda w^*\|^2 &= (\Phi^T(Y - \Phi \Phi^T \alpha^*) - \lambda \Phi^T \alpha^*)^T (\Phi^T(Y - \Phi \Phi^T \alpha^*) - \lambda \Phi^T \alpha^*) \\ &= (Y - \Phi \Phi^T \alpha^* - \lambda \alpha^*)^T \Phi \Phi^T (Y - \Phi \Phi^T \alpha^* - \lambda \alpha^*) = 0. \end{aligned}$$

Protože norma vektoru je nula právě tehdy, když je vektor nulový, dostáváme výsledek □

Gramova matice

Definujeme **Gramovu matici** (angl. **Gram matrix**) vztahem $\mathbf{G} = \Phi\Phi^T \in \mathbb{R}^{N,N}$.

Gramova matice je zjevně **symetrická** a protože

$$\mathbf{a}^T \mathbf{G} \mathbf{a} = \mathbf{a}^T \Phi\Phi^T \mathbf{a} = (\Phi^T \mathbf{a})^T (\Phi^T \mathbf{a}) = \|\Phi^T \mathbf{a}\|^2 \geq 0$$

pro každé $\mathbf{a} \in \mathbb{R}^N$, je také **pozitivně semidefinitní**.

$\text{RSS}_\lambda(\alpha)$ můžeme vyjádřit jako

$$\text{RSS}_\lambda(\alpha) = \|\mathbf{Y} - \mathbf{G}\alpha\|^2 + \lambda\alpha^T \mathbf{G}\alpha.$$

Definujme dále **jádrovou funkci** (angl. **kernel function**) $k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ jako

$$k(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^T \varphi(\mathbf{y}) \quad \text{pro každé } \mathbf{x}, \mathbf{y} \in \mathbb{R}^p.$$

Pro i, j -tou složku Gramovy matice platí

$$G_{i,j} = (\Phi\Phi^T)_{i,j} = \sum_{\ell=1}^M \Phi_{i,\ell} \Phi_{j,\ell} = \sum_{\ell=1}^M \varphi_\ell(\mathbf{x}_i) \varphi_\ell(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j),$$

tj. Gramova matice je **plně určena jádrovou funkcí**.

Predikce v duální reprezentaci

Předpokládejme, že jsme našli $\hat{\alpha}$ minimalizující $\text{RSS}_{\lambda}(\alpha)$.

Odpovídající \hat{w} minimalizující $\text{RSS}_{\lambda}(w)$ je dáno vztahem

$$\hat{w} = \Phi^T \hat{\alpha}.$$

Pro predikci Y v bodě x máme

$$\hat{Y} = \hat{w}^T \varphi(x) = \hat{\alpha}^T \Phi \varphi(x) = \sum_{i=1}^N \sum_{j=1}^M \hat{\alpha}_i \Phi_{i,j} \varphi_j(x) = \sum_{i=1}^N \hat{\alpha}_i k(x_i, x),$$

Tudíž nejenom účelovou funkci¹ $\text{RSS}_{\lambda}(\alpha)$ ale také predikce \hat{Y} můžeme vytvořit s využitím pouze jádrové funkce k .

Jak ukážeme, platí to i pro explicitní vyjádření $\hat{\alpha}$.

¹Název pro funkci, kterou minimalizujeme.

Explicitní řešení minimalizace

Najdeme **explicitní řešení** duální minimalizační úlohy.

Pro gradient $\nabla \text{RSS}_\lambda(\alpha)$ platí

$$\nabla \text{RSS}_\lambda(\alpha) = -2\mathbf{G}(\mathbf{Y} - \mathbf{G}\alpha) + 2\lambda\mathbf{G}\alpha.$$

Položíme-li ho roven nule, dostaneme normální rovnici

$$\mathbf{G}(\mathbf{Y} - \mathbf{G}\alpha - \lambda\alpha) = \mathbf{0}.$$

Protože je Gramova matice pozitivně semidefinitní, je matice $(\mathbf{G} + \lambda\mathbf{I})$ regulární a dostáváme

$$\hat{\alpha} = (\mathbf{G} + \lambda\mathbf{I})^{-1}\mathbf{Y}.$$

Tedy i optimální $\hat{\alpha}$ můžeme spočítat pouze na základě Gramovy matice a tedy jádrové funkce.

Shrnutí duální reprezentace

- Při trénování minimalizujeme $\text{RSS}_\lambda(\alpha)$ jakožto duální verzi reziduálního součtu čtverců danou vztahem

$$\text{RSS}_\lambda(\alpha) = \|\mathbf{Y} - \mathbf{G}\alpha\|^2 + \lambda\alpha^T \mathbf{G}\alpha,$$

kde $G_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$.

- Řešení **minimalizace** je pro každé $\lambda > 0$ rovno

$$\hat{\alpha} = (\mathbf{G} + \lambda\mathbf{I})^{-1}\mathbf{Y}.$$

- **Predikce** \hat{Y} v bodě \mathbf{x} je potom

$$\hat{Y} = \sum_{i=1}^N \hat{\alpha}_i k(\mathbf{x}_i, \mathbf{x}).$$

- Zároveň vidíme, že predikce \hat{Y} v bodě \mathbf{x} je vlastně váženou lineární kombinací hodnot bodů z trénovací množiny, kdy váhy jsou spočteny pomocí jádrové funkce. To umožňuje **intepretovat** výsledky predikce.

Jádrový trik

- Vytvořili jsme **ekvivalentní duální reprezentaci** celého modelu včetně účelové funkce pro trénování.
- V této reprezentaci se body $\mathbf{x}, \mathbf{x}_1, \dots$ z originálního příznakového prostoru \mathcal{X} vyskytují pouze ve tvaru skalárních součinů jejích transformací pomocí báзовých funkcí (tj. skalárních součinů v novém příznakovém prostoru), $\varphi(\mathbf{x})^T \varphi(\mathbf{y})$, které můžeme kompletně vyjádřit pomocí jádrové funkce $k(\mathbf{x}, \mathbf{y})$.
- Toto nahrazení skalárních součinů pomocí jádrové funkce se nazývá **jádrový trik** (angl. **kernel trick**).
- Přirozené rozšíření je rovnou **začít s jádrovou funkcí** bez explicitního zavádění báзовých funkcí.
- Tento přístup umožňuje implicitně pracovat v příznakových prostorech vysoké (nekonečné) dimenze.
- Z pohledu výpočetní náročnosti je dobré si uvědomit, že maticová inverze potřebná pro odhad \mathbf{w} resp. α má složitost $\mathcal{O}(M^3)$ resp. $\mathcal{O}(N^3)$.

Příklad využití jádrového triku

- Uvažujme například úlohu regrese nad 1000 obrázky o rozměrech $32 \times 32 = 1024$ pixelů ve stupních šedi.
- Při použití kvadratického jádra $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^2$ je ekvivalentní lineárnímu modelu bazových funkcí v 525 825 rozměrném prostoru.
- Tento počet odpovídajících bazových funkcí lze snadno určit, pokud položíme $x_0 = y_0 = 1$ a všimneme si, že

$$k(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n x_i y_i + 1 \right)^2 = \left(\sum_{i=0}^n x_i y_i \right)^2 = \sum_{i=0}^n \sum_{j=0}^n (x_i x_j) \cdot (y_i y_j)$$

což má $(n+1) \cdot (n+1)/2 = 525\,825$ rozdílných složek.

- Při použití jádrového triku a duální reprezentace musíme invertovat matici 1000×1000 namísto matice $525\,825 \times 525\,825$.

Lineární a polynomiální jádrové funkce

Nyní si představme **nejdůležitější příklady** používaných jádrových funkcí.

Když $\mathcal{X} = \mathbb{R}^p$ a $\varphi(x) = x$ pro všechna $x \in \mathbb{R}^p$ dostáváme **lineární jádro** (angl. **linear kernel**)

$$k(x, y) = x^T y.$$

Zobecněním je (nehomogenní) **polynomiální jádro** (angl. **polynomial kernel**)

$$k(x, y) = (x^T y + 1)^n.$$

Bázové funkce, které toto jádro **implicitně definuje** si ukažme na příkladu $p = 2$ a $n = 2$:

$$\begin{aligned}(x^T y + 1)^2 &= (x_1 y_1 + x_2 y_2 + 1)^2 \\&= (x_1 y_1)^2 + (x_2 y_2)^2 + 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 y_1 x_2 y_2 \\&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2) (1, \sqrt{2}y_1, \sqrt{2}y_2, y_1^2, y_2^2, \sqrt{2}y_1 y_2)^T.\end{aligned}$$

Což odpovídá 6 bazovým funkcím.

Gaussové jádro

Jednou z nejpoužívanějších jádrových funkcí je **Gaussové jádro** (angl. **Gaussian kernel** nebo **squared exponential kernel**), které je definováno jako

$$k(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}.$$

pro každé $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$.

Tato jádrová funkce je příkladem **radiální bákové funkce** nebo také **RBF jádra**, což znamená, že je to funkce pouze $\|\mathbf{x} - \mathbf{y}\|$.

Bývá ale také často označována pouze jako **RBF jádro** (viz např. ve [scikit-learn](#)).

Existuje i obecnější formulace Gaussového jádra, kde mohou být různé škály ve směrech různých příznaků.

Existují i další jádrové funkce, z nichž některé umí pracovat i s nečíselnými vstupy (např. s dvojicemi konečných řetězců).

Obecně, aby vše „dobře fungovalo“, musí být jádrová funkce **pozitivně semidefinitní symetrická funkce**, typicky také nezáporná. Co to přesně znamená si zde ovšem nebudeme podrobněji vysvětlovat.

Jádrové modely

- Uvažujme, že skutečný model, ze kterého pochází vysvětlovaná proměnná Y v bodě \mathbf{x} je

$$Y = f(\mathbf{x}) + \varepsilon,$$

kde f je nějaká neznámá funkce a ε je náhodná veličina s $E\varepsilon = 0$.

- V lineárním modelu báзовých funkcí hledáme f ve tvaru

$$f(\mathbf{x}) = w_1\varphi_1(\mathbf{x}) + \dots w_M\varphi_M(\mathbf{x}).$$

- V obecném **jádrovém modelu** (angl. **kernel machine**) s jádrovou funkcí k máme f ve tvaru

$$f(\mathbf{x}) = \sum_{j=1}^K \alpha_j k(\boldsymbol{\mu}_j, \mathbf{x})$$

kde $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \in \mathcal{X}$ jsou nějaké středové body.

- Speciální případ jádrových modelů (angl. **vector machines**) nastává, když středové body jsou body trénovací množiny, tj.

$$f(\mathbf{x}) = \sum_{j=1}^N \alpha_j k(\mathbf{x}_j, \mathbf{x}).$$

- Takový regresní model jsme dnes přesně získali pomocí jádrového triku.