

Mudel on ...

Mudel on empiiriliste andmete üldistus ja/või teoreetiline konstruktsioon, mis võimaldab prognoosida muutujate väärtusi.

- Mudel **üldistab ja lihtsustab**, ebaoluline jäetakse kõrvale (**üldistav funktsioon**).
- Mudeli koostamine annab teadmisi protsesside iseloomu ja statistiliste seoste kohta (**kirjeldav funktsioon**).
- Mudel annab näidiseid tegelikkuse **võrdlemiseks** teoreetiliste variantidega (juhu kui ...) (**võrdlusfunktsioon**).
- Mudel annab (formaalselt) põhjendatud **ennustuse (prognoosiv funktsioon)**.
- Mudelist saadud prognooside statistilist olulisust on võimalik hinnata, kui uurida lähteandmete ja korduvate prognooside varieeruvust (**tõestav funktsioon**).

Kas statistiline modelleerimine on kirjeldav või tõestav andmeanalüüs?

Säästvusreegel Parsimoonia ladina keeles **kokkuhoidlikkus**, parsimoonne = kokkuhoidlik, säästlik

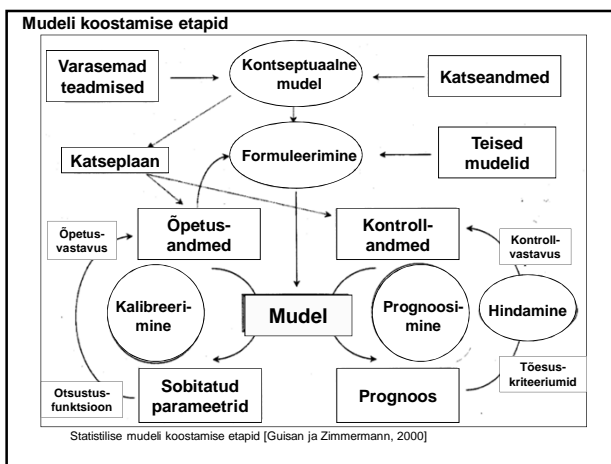
Parsimoonia reegel ehk Occami (või Ockhami) habemenuga (*Occam's razor*) — keskaegsele filosoofile Ockhami Williamile omistatav loogikaprintsiip:

Entia non sint multiplicanda praeter necessitatem
Üheski Occami kirjutises otseselt sellist sõnastus ei ole (Wikipedia).

keerukust ei tohiks eeldada ilma tungiva vajaduseta
ehk
kui mingit nähtust on võimalik seletada lihtsalt, siis ei ole põhjust kasutada keerukat seletust
ehk
tuleb eelistada lihtsamat seletust
ehk
tuleb olla nii laisk, kui võimalik.

Kuidas mõõta keerukust?
? Seletuses (mudel) kasutatavate tunnuste arv.
? Teisenduste arv.
? Mudeli koostamiseks kulunud tööaeg ja töötasu.

Kuidas mõõta seletuse õigsust?
? Vastavus empiirilise kogemusega.
? Vastavus varasemate teadmistega (ettekujutustega).



Staatilised statistilised mudelid

Pidevad mudelid (funktsioontunnus on pidev)

Lineaarsed mudelid

- Lihtsad lineaarsed mudelid
- Üldised lineaarsed mudelid
- Üldistatud lineaarsed mudelid
- Teised lineaarsed aditiivsed mudelid

Mittelineaarsed pidevad mudelid

Aegridade mudelid

Klassifikatsioonimeetodid (f-tunnus on nominaalne)

Klassifikatsioonipuud, regressioonipuud, hägusad (töenäosuslikud) klassifikatsioonid, diskriminantanalüüs, näidistele tuginevad klassifikatsioonid

Intellektitehnika (hinnangute aluseks on iteratiivne isekohanduv protsess — tehisope, võimalik on andmete pidev lisamine)

Tehisnärvivõrgud
Ekspertsüsteemid
Näidistele tuginevad süsteemid
Otsuste puud

Lihtsad lineaarsed meetodid ja mudelid

Funktsioontunnus on pidev

<p>Regressioonanalüüs pidevad argumenttunnused</p> <p>Graafiline kujutis: Regressioonijoon — joon läbi korrelatsioonivälja, millest hälvete (regressioonijääkide) summa on minimaalne. Regressioonijääkide mõõtmiseks on erinevaid variante.</p> <p>Vastavuspind (<i>response surface</i>) — Kahe ja enama pideva argumenttunnuse korral.</p>	<p>Kovariantsanalüüs (<i>analysis of covariance — ANCOVA</i>)</p> <p>Osa argumenttunnuseid on pidevad ja osa nominaalsed. Annab teavet, kuidas sõltub pideva argumenttunnuse prognoosimisvõime nominaalse argumenttunnuse klassidest.</p> <p>Graafiline kujutis: eraldi regressioonijoon kategoorilise muutuja iga väärtusklassi puhul.</p>
---	---

Dispersioonanalüüs (*analysis of variance — ANOVA*)
funktsioontunnus on pidev, kõik argumenttunnused **nominaalsed**.

Dispersioonanalüüs on meetod pideva tunnuse varieeruvuse jaotamiseks etteantud gruppide vaheliseks ja gruppide siseseks varieeruvuseks.

Lihtne lineaarne regressioon võrrandina Funktsioontunnuse ootus (keskväärtus).

Mudeli võrrand:
 $Y = b_0 + b_1 X_1 + \text{viga} \{ \sim N(0, \sigma^2) \}$ ehk $EY = b_0 + b_1 X$

Funktsioontunnus ehk sõltuv tunnus ehk prognoositav tunnus. Selle väärtusi arvutatakse.

Vabaliige (*intercept*) (funktsioontunnuse prognoos juhul, kui argumenttunnus=0).

Regressioonikordaja (*regression coefficient*). Näitab funktsioontunnuse muutust argumenttunnuse ühikulise muutuse korral. Sõltub tunnuste mõõtühikutest!

Normaaljaotusega juhuslik viga (saab kasutada prognoosile usalduspiiride lisamiseks). Muude faktorite mõju

Argumenttunnus ehk seletav tunnus (selle väärtused on ette antud).

Korrelatsioon on vastastikune seos.
Regressioon on ühe tunnuse sõltuvus teisest, s.t. et reeglina ei ole regressioonimudel pööratav.

Regressioonimudeli näide

$$S = -5.26 + 6.36 \cdot \text{elupymbr} + 0.035 \cdot \text{teedzn} + 0.013 \cdot \text{kaezn} + 0.047 \cdot \text{asustuszn}$$

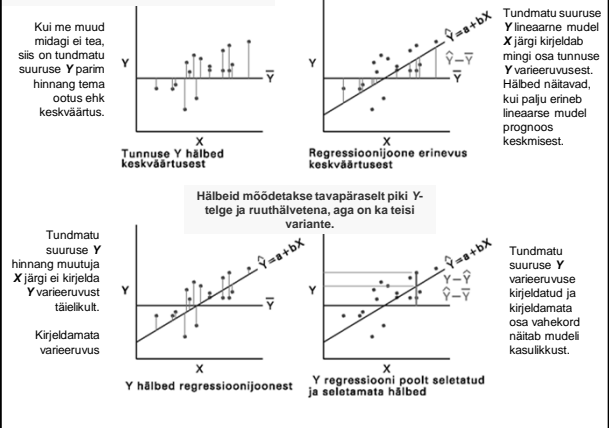
$S = \text{logit}(\text{põdra esinemistõenäosus}) = \text{koha sobivus põdra elupaigana}$

elupymbr – põdra elupaiga kaugusega kaalutud osakaal ümbruses,
teedzn – kaugus maanteedest,
kaezn – kaugus kaevanduste lõhketöödest,
asustuszn – kaugus inimasustusest.

$$P = \exp(S) / \{1 + \exp(S)\}$$

NB! Võrrandis kasutatud tähistus tuleb ära seletada!

Lineaarse regressiooni hälbed



Andmekogumi varieeruvuse kirjeldamiseks arvutatakse **dispersioon** keskväärte suhtes.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}$$

Keskväärte kui juhusliku muutuja esmane esindaja (ootus)

Standardviga

Mudeli puhul leitakse prognoosi hälbed vaatlustest.

Kui soovitakse hinnata mudeli täpsust **väljaspool õpetusandmeid** (väljaspool valimit), peaks dispersiooni ja standardhälbe arvutamisel nimetajas olema vabadusastmete arv = vaatluste arv (n) – tunnuste arv mudelis (t).

Dispersiooni arvutamisel valimis jagati (n – 1)

$$s^2 = \frac{\sum_{i=1}^n (x_i - \text{pred}_i)^2}{n - t}$$

$$RMSE = \sqrt{s^2}$$

Ruutkeskmine viga (*standard error of residuals = standard residuals = root mean squared error = RMSE*) ehk prognoosijääkide standardhälve ehk **standardviga** kirjeldab prognoosi oodatavat täpsust.

Vaatluste varieeruvust kirjeldab standardhälve.

Vabadusastmete arv

Sõltumatute muutujate arv mudelis või tunnuruumi dimensionide arv.

- 0 mõõtmelises ruumis (punkt) koha määratust ei ole
- 1 mõõtmelises ruumis (joonel) on koht määratav ühe arvuga (kaugus 0-punktist)
- 2 mõõtmelises ruumis (tasapinnal) on koha määratamiseks vaja kahte parameetrit (x ja y või suund ja kaugus või kaks suunda)
- 3 mõõtmelises ruumis on koha määratamiseks vaja kolme parameetrit (x, y ja z või kaks suunda ja kaugus või suund ja kaks kaugust või kolm suunda)

Lineaarne mudel: $Y = a + bx$
 ühe seletava tunnusega (sobitavad parameetrid a ja b) kirjeldab kahe andmepunkti paiknemist hälveteta. Läbi kahe punkti saab alati tõmmata sirge.

Lineaarne mudel: $Y = a + b_1x_1 + b_2x_2$
 kahe seletava tunnusega (x_1 ja x_2 , mudeli parameetrid a, b_1 ja b_2) kirjeldab kolme andmepunkti paiknemist hälveteta.

Kui vaatlusi ei ole rohkem kui mudelis vabadusastmeid, saab kogu varieeruvuse ära seletada, aga me ei tea midagi mudeli kehivusest väljaspool neid vaatlusi. => **parsimoonia reegel (eelistati lihtsast seletust).**

Vähimruutude meetodil sobitatud lihtsate lineaarsete mudelite eeldused

- Faktorite mõjud ja koosmõjud on aditiivsed (liituvad), faktorite vahel puuduvad teistsugused koosmõjud.
- Kõik olulised seletavad tunnused on mudelis sees.
- Mudeli viga on konstantse dispersiooniga kõigi argumenttunnuse väärtuste korral (ei ole heteroskedastsust).
- Vaatlused on sõltumatud ja ühekordsed.** Vaatlusvead on sõltumatud (ei sõltu vaatluste sarnasusest).
- Mõõtmisvigade keskväärte on nihketa (valimi keskväärte on üldkogumi keskväärte parim hinnang, valim on **esinduslik**).
- Vead on mudeli suhtes normaajaotusega.
- Seos funktsioon- ja argumenttunnuste vahel on lineaarne või on andmed lineaarsuse saavutamiseks teisendatud.

Kui mõni eeldus ei ole täidetud, siis peaks otsima teisi modelleerimise meetodeid

Mittelineaarsed regressioonimudelid

Näiliselt mittelineaarsed (loomult lineaarsed) mudelid (<i>intrinsically linear</i>)	Funktsioonitunnuse teisendamisega lineaarselt muutuvad mudelid	Päriselt mittelineaarsed mudelid
Argumenttunnuse teisendamisega lineaarselt muutuvad mudelid	Funktsioonitunnuse teisendamisega lineaarselt muutuvad mudelid	Mudeli viga ei allu muutujatega samale teisendusele
Polünoomregressioonid	$N = e^{b_1 \times a e^{b_2}}$	$N = e^{b_1 \times a e^{b_2}} + \text{viga}$
$N = b_0 + b_1x + b_2x^2$	$\log(N) = b_1 \times a e^{b_2}$	Murdepunktiga mudelid
$N = b_0 + b_1x_1 + b_2x_2$	$N = e^{b_1 \times a e^{b_2}} \times \text{viga}$	$N = b_0 + b_1x (x < 100)$
	$\log(N) = b_1 \times a e^{b_2} + \text{viga}$	$N = b_{00} + b_2x (x \geq 100)$
	Log-mudelid on vead aditiivsed	

Erinevus on eeldatavas vigade jaotuses

Üldised lineaarsed mudelid (GLM)

- **Mitu funktsioontunnust** — funktsioontunnuste **plaanimaatriks**
- Argumenttunnuste mõjud võivad olla omavahel **korreleerunud**
- Andmed võivad sisaldada **korduvmõõtmisi**
- **Pidevad ja nominaalsed** argumenttunnused võivad olla kombineeritud
- Eeldatakse funktsioontunnuse (regressioonivigade) **normaaljaotust**
- Funktsioontunnuste **teisendused** tuleb teha **enne** mudeli sobitamist

Nominaalsed argumenttunnused **kodeeritakse** kas numbriliselt või täheleiselt.

Standardiseeritud kodeerimise (*sigma restricted*) korral omistatakse binarisele tunnusele koodid 1 ja -1.

Üleparametriseeritud mudelis kodeeritakse kõik variandid omaette tunnusteks väärtustega 0 ja 1 (esineb / puudub).

Üldistatud lineaarsed mudelid (GLMZ)

- ⊗ **Mitu funktsioontunnust** — funktsioontunnuste **plaanimaatriks**.
- ⊗ Mõjud võivad olla omavahel **korreleerunud**.
- ⊗ Pidevad ja **nominaalsed** tunnused võivad olla kombineeritud.
- ⊗ Andmed võivad sisaldada **korduvmõõtmisi**.
- **Funktsioontunnuse jaotus võib olla mitmesugune**, k.a. nominaalne.
- Jaotust ja lineariseerivat teisendust (seosefunktsiooni) tuleb eelnevalt **teada**.

Seosefunktsioonid (link functions)

normaal-, gamma-, pöördnormaal- ja Poissoni jaotus:
 Identity link: $f(z) = z$ kui eeldatakse vigade normaaljaotust
 Log link: $f(z) = \log(z)$ loendusandmed, Poissoni jaotusega vead
 Power link: $f(z) = z^\alpha$, funktsioontunnus astendatakse

binoom- ja järjestatud multinomialjaotus:
 Logit link: $f(z) = \log(z/(1-z))$
 Probit link: $f(z) = \text{invnorm}(z)$ *invnorm* on standardiseeritud kumulatiivse normaaljaotuse pöördväärtus
 Log-log link: $f(z) = -\log(-\log(z))$
 Complementary log-log link: $f(z) = \log(-\log(1-z))$

multinomialjaotus:
 Generalized logit link: $f(z_1|z_2, \dots, z_c) = \log(x_1/(1-z_1-\dots-z_c))$ mudelis on $c+1$ kategooriat

Logitfunktsioon on logaritm šansside suhtest (odds ratio) ehk esinemistõenäosuse puudumistõenäosuse vahekorrasst ehk tõenäosuste suhtest. **Logitfunktsioon kirjeldab ühesuunalist tendentsi binarsetes andmetes, kus mõlema variandi esinemise tõenäosus on vahemikus 0...1.**

Probitregressiooni mudel on $y = N(b_0 + \sum b_i x_i)$, kus N on standardiseeritud kumulatiivne normaaljaotus (normaaljaotuse kõvera alune pind). Probitfunktsioon kasutab normaaljaotust pideva muutuja teisendamiseks vahemikku 0...1.

Logitmodel

Logitfunktsioon on logaritm šansside suhtest (odds ratio) ehk esinemise puudumise vahekorrasst ehk tõenäosuste suhtest. Logit-funktsioon kirjeldab ühesuunalist tendentsi binarsetes andmetes, milles mõlema variandi esinemise tõenäosus on vahemikus 0...1.

Esinemistõenäosuse

$$\text{logit} = \log \frac{p}{1-p}$$

põdra esinemistõenäosus

logitprognoos

Logitmodeliga prognoositud esinemistõenäosus

$$p = \frac{e^{\text{logit}(p)}}{1 + e^{\text{logit}(p)}}$$

Kas logitmodel on lineaarne mudel või ei ole?
 Sõltub vigade jaotusest. Kui logit-teisendus järel saadakse normaaljaotusega häbed, siis on mudel näiliselt mittelineaarne, aga sisuliselt siiski lineaarne.

Põdraga vaatlusalad
 Põdrata vaatlusalad

Probleemid lineaarsete mudelite kasutamisel

Probleem	Tagajärg	Kontrollimisviis	Abinõu
Jäägid ei ole normaaljaotusega.	F-statistikul baseeruvad testid ei anna õigeid tulemusi. Nihega hinnangud.	Jääkide graafik, normaaljaotuse testid.	Funktsioontunnuste teisendamine. Üldistatud lineaarse mudeli või mitteparameetriselise meetodi kasutamine.
Jääkide hajuvus on ebahõlne.	Regressiooni täpsuse hinnangud hälbevad.	Jääkide graafik.	Funktsioontunnuste teisendamine. Kaalude omistamine vaatlustele. Mitteparameetriselised meetodid.
Jäägid ei ole sõltumatud.	Regressiooni tugevus ülehinnatud.	Jääkide graafik, Morani test ruumiandmete puhul, Durbin-Watsoni test kõigi järjestuvate andmete puhul.	Iteratiivsed sobitamismeetodid. Üldistatud lineaarse mudeli.
Seose mittelineaarsus.	Mudeli vähene kirjeldav võime, sisutud tulemused.	Korrelatsiooniväärt.	Muutujate teisendamine. Mittelineaarne mudel. Mitteparameetriselised meetodid.
Argumenttunnuste multikollineaarsus.	Ebastabiilne mudel.	Korrelatsioonimaatriks.	Argumenttunnuste arvu vähendamine, argumenttunnuste teisendamine, rüüregressioon.
Võimalika argumenttunnuste suur hulk.	Raske valida optimaalset tunnuste komplekti.	Tunnuste loend.	Samm-sammuline regressioon, muutujate teisendamine, mudeli lihtsustamine.
Erindid.	Mudel on erinditest liigselt sõltuv.	Kirjeldav andmeanalüüs.	Erindite eemaldamine. Mitteparameetriselised meetodid.
Vead andmetes.	Sisutud tulemused.	Kirjeldav andmeanalüüs.	Vigaste andmete eemaldamine või asendamine.
Lünklikud andmed.	Olemasolevate andmete ebaefektiivne kasutus.	Risttabel.	Puuduvate andmetega vaatluste eemaldamine või asendamine prognoositud väärtustega või keskmisega.
Nominaalsed tunnused.	Normaaljaotust eeldav regressioon ei ole kasutatav.	Väärtuste jaotumise vaatlus.	Logistiline regressioon või muu üldistatud lineaarse mudeli variant.

Ebasobiva meetodi rakendamise tagajärjed

- ❖ Mitteparameetriselise meetodi kasutamine, kui parameetriselise meetodi eeldused on täidetud => otsitav seos või erinevus jääb avastamata.
- ❖ Kaalude omistamine vaatlustele => tulemuste (meelevaldne) muutmine.
- ❖ Keerukama meetodi kasutamine kui piisaks lihtsamast => uurimus ei jäta usaldusväärselt mujalt, risk eksida keeruka meetodi paljude parameetrite ja variantide hulgast sobivaima valikul, suurem arvutuste maht ja kulutatud aeg, tulemusi on keerukam kontrollida.
- ❖ Tunnuste teisendamine => tunnuste muutmispõrkkondade mõju tulemustesse muutub.
- ❖ Erindite eemaldamine => erindid ei pruugi olla vigased vaatlused, võib olla esindavad mingit olulist nähtust. Tulemuse sobitamine oma suva järgi erindeid valides.
- ❖ Vigaste andmete asendamine => subjektiivsus.
- ❖ Puuduvate väärtuste asendamine keskmistega => varieeruvus väheneb. Kas varieeruvust on tarvis hinnata?
- ❖ Puuduvate väärtuste asendamine juhuslike väärtustega => seoste, varieeruvuse ja keskmise muutmise.
- ❖ Lünklike vaatluste eemaldamine => kasutamata jäävad ka mõõdetud tunnuste väärtused.
- ❖ Väärtuste juhuslik segipaikamine (permuteerimine) => kasutatav vaid jaotuse genereerimiseks nullhüpoteesi kehtimisel.

Näidiste järgi prognoosimine

Näidiste järgi saab nii **klassifitseerida** (klassikuluuvust prognoosida) kui ka pideva muutuja väärtusi (puude kätvust, liigi esinemistõenäosust, mingi protsessi intensiivsust) **prognoosida**.

Kuna enamasti täpselt vastavat näidist ei leidu, siis kasutatakse pidevate andmete puhul prognoosina sarnasemate näidiste **sarnasusega kaalutud keskmist**.

Sarnasuse järgi prognoosimisel kasutatakse:

kõiki vähemalt teatud sarnasusega näidiseid **teatud hulka kõige sarnasemaid näidiseid**

d-lähima naabri meetod ehk d-NN meetod **k-lähima naabri meetod ehk k-NN meetod**

d-lähima naabri meetodi puhul varieerub näidiste hulk, mille järgi prognoos arvutatakse. Kui näidised ei paikne tunnusruumis ühtlaselt, on hinnangu usaldusväärsus varieeruv.

Kui piisavalt sarnaseid näidiseid ei ole, siis prognoosi ei saa.

k-lähima naabri meetod on kasutatav igasuguse andmete tiheduse juures. k-lähimat naabrit on lähimad naabrid ka siis, kui nad väga kaugel on.

Prognoos saadakse ka juhul, kui sarnaseid näidiseid ei ole, hinnangu täpsus varieerub.

Näidiste kaalumine

Klassifitseerimine näidiseid kaalumata

Kategoriline muutuja

Ühe kõige sarnasema näidise järgi tuleb vaatlus lugeda klassi koodiga 5, nelja näidise korral klassi 2.

Pideva tunnuse prognoos kaalutud näidiste järgi

Pidev muutuja

Joonisel olevate kaaludega (w) kaalutud näidise korral tuleks vaatlusele omistada väärtus 3.5. Kaalud on määratud kaugusega tunnusruumis.

Kas siin kujutatu on k-NN või d-NN meetod?

Näidiste baas

Näidiste järgi prognoosimiseks on mudeli asemel tarvis **näidiste baasi** ja **tarkvara**, mis leiaks andmebaasist kiiresti iga juhtumi puhul soovivaimad näidised ning tuletaks näidistest hinnangu või ennustuse.

Näidistega sarnasusele tuginevaks tehisepeeks, hinnangute ja hinnangukaartide arvutamiseks on geoinformaatika õppetoolis loodud tarkvarasüsteem *Constud*.

<http://www.geo.ut.ee/CONSTUD>

Aja modelleerimine

Aeg kui **funktsioontunnus** uuritavaks suuruseks on aeg, mis kulub teatava sündmuseni (surmani, haigestumiseni, detaili riknemiseni jmt). Näiteks **(elu)kestusanalüüs**.

Aeg kui **müratunnus** Vaatlused peaksid pärinema samast ajamomendist, kuid mingitel põhjustel on vaatlusaeg veninud pikemaks. Üldine probleem kaartide ja fotode kasutamisel ja ka väliandmete kogumisel. Lahendus võib olla aja lisamine argumenttunnuste hulka.

Aeg kui **argumenttunnus** Ajas muutuva tunnuse või tunnuste väärtusi nimetatakse **aegreaks**. Aegriidade mudelites on aeg põhiline (tihti ainus) argumenttunnus.

Tavaliselt on aegreas: iga ajahetke kohta vaid üks vaatlus; vaatluste koguarv suur; aega käsitletakse ühtlaselt diskreetsena (fikseeritud sammuga aegriida).

Aegrea komponendid

Andmete varieeruvus aegreas → **Argumenttunnustega seotud varieeruvus.**
Seda soovitakse mürast eraldada

Müra (juhuslikud hälbed)
Prognoosida saab müra keskmist tugevust ja üksikvaatluse ootust, aga mitte üksikvaatluste juhuslikke hälbeid.

Aegriidade modelleerimiseks kasutatakse:
trendi (kindlasuunalist muutust),
sesoonsust (tsüklilisust),
autokorrelatsiooni (endaga sarnasust).

Aegriida on siin joonisel komponentide summa

Mudeli kontrollimine ja parima mudeli valik

Otsustusfunktsioon (objective function) — parema mudeli valimise aluseks olev kriteerium.

Liigsobitumine (overfitting) — olukord, mileni jõutakse mudeli liigsel sobitamisel andmetele. Liigsobitunud mudeli lahend annab õpetusandmete puhul tunduvalt parema prognoosi kui väljaspool õpetusandmeid.

Treeningandmestik ehk **õpetusandmestik (training data set)** — vaatlusandmed, mille järgi mudelit kalibreeritakse.

Kontrollandmestik (validation data set) — vaatlusandmed, mille järgi mudeli kalibreerimist ja seletavat võimet kontrollitakse.

Ristkontroll (cross-validation) — prognoosi täpsuse kontroll õpetusandmetest **sõltumatu** kontrollandmestikuga.

Kontrollandmestikele võib olla üks või mitu.

Vaatluste ükshaaval kõrvalejätmine (LOOC — leave-one-out cross-validation).

Üleõppimine ja tühitargutamine tekib siis, kui vähestest faktidest üritatakse saada kauguleulatavaid järeldusi.

Prognoositäpsuse hindamine ja mudelite võrdlemine

Möödupuud (otsustusfunktsioon)

Normaaljaotusega muutuja — ruutkeskmine viga või determinatsioonikordaja, s.o. mudeli poolt ära kirjeldatud varieeruvuse osa muutuja varieeruvusest keskvaartuse suhtes. Sobib pideva tunnuse mudelitel puhul, kui funktsioontunnuse hälvete ruutude kasutamine on põhjendatud.

Muudel juhtudel:
Hinnatakse jaotuse (klassifikatsioonide) vastavust —
Kapa, vigade maatriks, χ^2 test, samasuskordajad, jt. meetodid (tõepärasuhte, tõepärasuhte, Akaike informatsiooni kriteerium).

Nendest oli juttu eelmistes loengutes

Möötmismeetodid:

- 1) ristkontroll
- 2) üksikvaatluste kordamööda kõrvalejätmine (*jack-knife*),
- 3) korduvad tagasipandad juhuslikud valimid (*bootstrap*),
- 4) juhuandmetest saadud prognooside varieeruvusega võrdlemine.

Tulemus

Treeningtäpsus (training accuracy)
Hinnatakse samadest andmetest, mille järgi mudel koostati.

Kontrolltäpsus (test accuracy)
Hinnatakse treeningandmetest sõltumatul kontrollandmestikust.

Võiks meelde jääda

- Mudel on **teooriale** tuginev üldistus. Mudel ei teki vaid andmetest.
- Säästvusreegli kohaselt tuleb eelistada lihtsamat mudelit ja seletust.
- Kalibreerimisel saadakse andmetele sobitatud mudeli parameetrid.
- Funktsioontunnuse ootust prognoosiv lihtne lineaarne regressioonivõrrand sisaldab
 - funktsioontunnust,
 - vabaliiget,
 - regressioonikordajaid,
 - argumenttunnuseid.
- Regressioonivõrrandi pööramisel ei saada oodatavalt parimat lahendit.
- Logit on logaritm esinemise ja puudumise tõenäosuste suhtest.
- Sarnasuse järgi prognoosimisel mudelit ei looda.
- Sõltumatu kontrollandmestik võimaldab hinnata mudeli üldkehtivust (kehtivust väljaspool õpetusandmeid).